Logarithm with correct rounding

université de BORDEAUX

Master Thesis

Master in Sciences and Technologies, Specialty in Mathematics, Cryptology and Computer Security.

Author

Sidali Zitouni-Terki <sid-ali.zitouni-terki@etu.u-bordeaux.fr>

Supervisor

Paul Zimmermann < Paul . Zimmermann@inria.fr>

Tutor

Gilles Zémor <gilles.zemor@u-bordeaux.fr>

september 2nd 2022

Declaration of authorship of the document

I hereby certify that this material, which I now submit for assessment on the programme of study leading to the award of the Master in *Sciences and Technologies*, Specialty in *Mathematics* or *Computer Science*, Track *Cryptology and Computer Security*, is entirely my own work, that I have exercised reasonable care to ensure that the work is original, and does not to the best of my knowledge breach any law of copyright, and has not been taken from the work of others save and to the extent that such work has been cited and acknowledged within the text of my work.

Date and Signature

Remerciement

Remerciements particuliers à mon superviseur Paul Zimmermann pour son aide et son soutien. Je le remercie également pour sa patience à mon égard. Je remercie mon tuteur Gilles Zémor et tous les autres enseignants pour leurs apports théoriques et méthodologiques et leurs partages de connaissances qui m'ont fait bénéficié tout au long de mon parcours universitaire. Je tiens également à exprimer mes sincères remerciements et ma gratitude aux membres du jury pour leur lectures et leur évaluation du travail.

Je remercie tout particulièrement :

Mes chers parents pour leur motivation et leur soutien constant et à qui je souhaite une longue vie pleine de santé.

Ma chère femme et ma belle-famille dont les mots d'encouragement m'ont motivé à reprendre et continuer mes études.

Mes filles bien-aimées, Farah et Tesnime.

Mes sœurs et mon frère qui m'ont soutenu dans ce travail.

Mes amis avec qui j'ai partagé des souvenirs inoubliables.

Résumé

L'implémentation des fonctions mathématiques avec arrondi correct est un sujet important en arithmétique flottant. Accéder à de telles fonctions nous permettrais d'avoir plus de précisions sur les calculs avec des nombres flottants. Nous avons des algorithmes, comme Fast2Sum et DEKKER-PRODUCT, qui nous aide à avoir plus de précisions. Grâce à ces algorithmes, nous pourrions implémenter le logarithme avec arrondi correct et avoir une fonction qui approche aux mieux au résultat du logarithme. Ce rapport nous montre d'abord: comment on peut utiliser les quatre modes d'arrondis qui sont l'arrondi au plus proche, arrondi au plus proche de zéro, arrondi au plus proche de l'infini et l'arrondi au plus proche de moins l'infini, et comment a été calculé la précision de chaque algorithme utilisé pour le logarithme. Ensuite, nous expliquons comment cette fonction de mathématique a été implanté. Ce logarithme avec arrondi correct a réussi le test sur le millions de pires cas.

Abstract

The implementation of mathematical functions with rounding correct is an important topic in floating point arithmetic. Access to such functions would allow us to have more precision on the calculations with floating numbers. We have algorithms, like **Fast2Sum** and **DEKKER-PRODUCT**, which help us to have more precisions. Due to these algorithms, we could implement the logarithm with rounding correct and have a function that best approximate the result of the logarithm. This report shows us first: how one can use the four modes roundings which are rounding to the nearest, rounding to the nearest to zero, rounding to nearest infinity and rounding to nearest minus infinity, and how was calculated the accuracy of each algorithm used for the logarithm. Inaddition, it explains how this function of mathematics were implanted. This logarithm with correct rounding passed the test on the millions of worst cases.

Contents

| C | onter | nts | 6 | | | | |
|----|-----------------|---|------------|--|--|--|--|
| In | \mathbf{trod} | uction | 9 | | | | |
| 1 | Floa | ating point arithmetic | 11 | | | | |
| | 1.1 | Definitions | 11 | | | | |
| | 1.2 | Necessary tools | 15 | | | | |
| | 1.3 | Notations | 18 | | | | |
| 2 | Ste | ps to calculate crlog | 19 | | | | |
| | 2.1 | Special cases | 19 | | | | |
| | 2.2 | Argument reduction | 19 | | | | |
| | 2.3 | polynomial approximation and evaluation | 24 | | | | |
| 3 | Оре | erators on Double-Double numbers | 27 | | | | |
| | 3.1 | Addition Operators | 27 | | | | |
| | 3.2 | Multiplication Operators | 40 | | | | |
| 4 | Оре | Operators on Triple-Double numbers | | | | | |
| | 4.1 | Addition Operators | 45 | | | | |
| | 4.2 | Multiplication Operators | 57 | | | | |
| 5 | \log | | 77 | | | | |
| | 5.1 | crlogfast | 77 | | | | |
| | 5.2 | crlogaccurate | 77 | | | | |
| | 5.3 | crlogadvanced | 77 | | | | |
| | 5.4 | log | 77 | | | | |
| C | onclu | ısion | 7 9 | | | | |

| CONTENTS | 7 |
|--------------|-----|
| Annexes | 83 |
| Bibliography | 123 |

Introduction

. This thesis aims to explain in a mathematical way the steps for the implementation of the logarithm with correct rounding on floating numbers. This function is one of the functions of the COREMATH 1 project. This project implements mathematical functions with correct rounding to be able to integrate into mathematical libraries for the new revision of the standard IEEE754.

The definition of correct rounding given by the IEEE754 standard is as follows "Given a mathematical function f and a floating point number x, the correct rounding of f(x) is the floating number y closest to f(x) according to the given rounding mode (nearest, towards zero, towards $-\infty$ or towards $+\infty$)". This standard imposes the correct rounding for the four elementary arithmetic operations which are addition, subtraction, multiplication and division. But it does not impose for mathematical functions. For now, there is no mathematical library that gives us exactly the correct rounding.

This project already has the implementation of this logarithm for single precision (binary32 format of IEEE754). The calculation steps of our logarithm will be in function of double precision (binary64 format of IEEE754). We will use basic floating point algorithms which are **FastSum** and **DEKKER-PRODUCT**.

My research paper will be devited into five chapters.

In the first chapter, we explain floating point and the IEEE754 standard with some definitions. Then we integrate some arithmetic tools that will be used for calculations in chapters 4, and 5. Then we give some notation rules.

In the second chapter, we detail the steps of the $cr \log$ first, we talk about the special cases. Then, we explain the argument reduction with the algorithms of **Tang** and **Gal**. After, we calculate and evaluate the ap-

https://core-math.gitlabpages.inria.fr

10 CONTENTS

proximation polynomial thanks to the formula of **Taylor**.

For the third chapter, we define the addition and multiplication algorithms which will give us results in **Double-Double** type and calculate their relative errors for each of these functions.

In the fourth chapter, we define other addition and multiplication algorithms which give us results in **Triple-Double** type and we calculate their relative errors.

The last chapter explain how the logarithm will unfold with its three $cr \log$.

Chapter 1

Floating point arithmetic

The next defintion take from [14], [16] and [19].

1.1 Definitions

Definition 1: Floating point is a way of writing real numbers. This expression is mostly used for computers. It is different from the fixed point. It is represented with a precision, base(radix), sign, mantissa and exponent. Let x be a real number, so in this way we have $x = s.M.\beta^{e-p+1}$ with s the sign, M the mantissa, β a base(radix), e the exponent and p the precision.

p represents the number of bits of the **mantissa**. In our research, we prefer to use E=e-bias to be possible to have a negative exponent and also to have $E_{min}=1-E_{max}$.

Floating points are composed of several precisions given by the **IEEE 754** standard.

The next defintion take from [21].

Definition 2: IEEE 754 is a floating point standard created by the Institute of Electrical and Electronics Engineers. It is the most used standard for computer calculations. This standard imposes formats for each floating number with its representations of sign, mantissa, exponent and also the five rounding modes used.

IEEE 754-2008 defines three possible formats for base of 2 of floating point: Binary32, Binary64 and Binary128.

- Binary32 (single precision) (Figure 1.1):
 Binary32 is stored on 32 bits, with 1 bit for the sign, 8 bits for the exponent and 23 bits for the mantissa without its strong bit (24 bits).
- Binary64 (double precision) (Figure 1.2):
 Binary64 is stored on 64 bits, with 1 bit for the sign, 11 bits for the exponent and 52 bits for the mantissa without its strong bit (53 bits).
- Binary128 (Quadruple precision) (Figure 1.3):
 Binary128 is stored on 128 bits, with 1 bit for the sign, 15 bits for the exponent and 112 bits for the mantissa without its strong bit (113 bits).

In the 3 Schematics below of the 3 formats, we have the leftmost bit is the strong bit and the rightmost is the weakest bit:

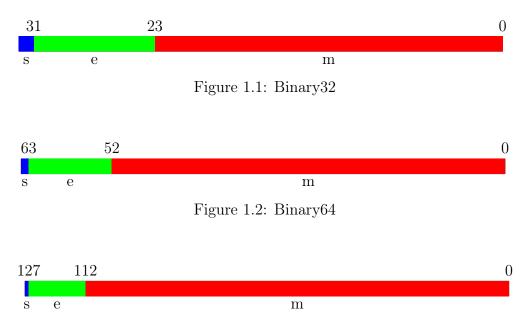


Figure 1.3: Binary128

We also have **Binary80** which is an **extended double precision** and two others possible formats for base of 10:

- decimal64 is stored on 64 bits.
- decimal 128 is stored on 128 bits.

In C programming, Binary32 represents Float, Binary64 is Double, Binary80 is long-double and Binary128 is float128.

Throughout our research, we will only use **Binary64**.

The real numbers are transformed into a floating number and we need to round them.

According to the IEEE 754 standard, there are four rounding modes among them, three are direct rounds:

- RNDN: Rounding to nearest.
 - It has two possibilities and the only difference between them is if the number falls in the midway.
 - **Ties to even**: it rounds to the nearest even floating number.
 - Ties to away: it rounds the nearest to the largest absolute value.
- **RNDZ**: Directed Rounding to 0 (truncation).
- RNDU: Directed Rounding to $+\infty$ (rounding up or ceiling).
- RNDD: Directed Rounding to $-\infty$ (rounding down or floor).

The examples below shows us the four rounding modes.

| Real Number | +140.8215064611465 | +665.5752955525412 |
|--------------------|--------------------|--------------------|
| RNDN(Ties to even) | +140.821506461146 | +665.575295552541 |
| RNDN(Ties to away) | +140.821506461147 | +665.575295552541 |
| RNDZ | +140.821506461146 | +665.575295552541 |
| RNDU | +140.821506461147 | +665.575295552542 |
| RNDD | +140.821506461146 | +665.575295552541 |

Table 1.1: roundings with positive reals

We have the special values:

- there exist two 0: +0 and -0
- Infinity: $-\infty$ and $+\infty$
- NaN (Not a Number): result of the invalid operations.(for example $:0/0, \sqrt{-7},...)$

| Real Number | -140.8215064611465 | -665.5752955525412 |
|--------------------|--------------------|--------------------|
| RNDN(Ties to even) | -140.821506461146 | -665.575295552541 |
| RNDN(Ties to away) | -140.821506461147 | -665.575295552541 |
| RNDZ | -140.821506461146 | -665.575295552541 |
| RNDU | -140.821506461146 | -665.575295552541 |
| RNDD | -140.821506461147 | -665.575295552542 |

Table 1.2: roundings with negative reals

- **Subnormal** Number: it is a number which has as exponent = 0 and a pseudo-mantissa. (difference between the Mantissa and the pseudo-mantissa those which does not have a hidden 1)
- Overflow: according to [1] "This exception is caused when the result is too large in absolute value to be represented".
- Underflow: according to [1] "This exception is caused when the actual result is too small in absolute value to be represented in the chosen format".

Definition 3: ulp(x)

According to **William Kahan** (1960) in [19] The original definition of ulp() by "ulp(x) is the gap between the two floating-point numbers nearest to x, even if x is one of them".

ulp(x) is the acronym of Unit of the Last Place.

In the following we will use Goldberg's definition:

Let x a **floating-point** number, $x = d_0d_1d_2d_3d_4...d_{p-1}\beta^e$. we therefore have an error to represent it : $|d_0d_1d_2d_3d_4...d_{p-1} - \frac{x}{\beta^e}|$ which is the unit of the place.

Corollary 1 If x is a floating-point number, so we have :

$$ulp(x) \le 2^{-52}|x|$$

.

1.2 Necessary tools

1.2.1 properties

The properties are from [19].

Property 1 If $X = RNDN(x) \Rightarrow |X - x| \le \frac{1}{2}ulp(x)$

Property 2 If $X = RNDN(x) \Rightarrow |X - x| \le \frac{1}{2}ulp(X)$

Property 3 If $X \in \{RNDD(x), RNDU(x), RNDN(x)\} \Rightarrow |X - x| \leq ulp(x)$

Property 4 If $X \in \{RNDD(x), RNDU(x), RNDN(x)\} \Rightarrow |X - x| \leq ulp(X)$

Property 5 If x > 0 then RNDZ(x) = RNDD(x). If x < 0 then RNDZ(x) = RNDU(x).

We need to calculate the relative error to know the accuracy of our algorithms.

Corollary 2 If $|X - x| \le \alpha.ulp(x)$ with $\alpha \in \mathbf{R} \Rightarrow \frac{|X - x|}{|x|} \le \alpha.2^{-52} \Rightarrow |X| \le |x|(1 + \alpha.2^{-52}).$

PROOF We suppose $|X - x| \le \alpha.ulp(x)$, as $ulp(x) \le 2^{-52}.|x|$ so we have:

$$|X - x| \le \alpha \cdot 2^{-52} \cdot |x|$$

$$\frac{|X-x|}{|x|} \le \alpha.2^{-52}$$

According to the triangle inequality, we have that $|X - x| \ge |X| - |x| \Rightarrow$

$$|X| - |x| \le \alpha .2^{-52}.|x|$$

$$|X| \le \alpha . 2^{-52} . |x| + |x|$$

$$|X| \le |x|(1 + \alpha \cdot 2^{-52}).$$

Also need to calculate with ulp(X)

Corollary 3 If $|X - x| \le \alpha.ulp(X)$ with $\alpha \in \mathbf{R} \Rightarrow |X| \le |x| + \alpha.ulp(X)$.

PROOF We suppose $|X - x| \le \alpha.ulp(x)$, as $ulp(X) \le 2^{-52}.|X|$ so we have:

$$|X - x| \le \alpha . 2^{-52}.|X|$$

$$\frac{|X-x|}{|X|} \le \alpha.2^{-52}$$

According to the triangle inequality, we have that $|X - x| \ge |X| - |x| \Rightarrow$

$$|X| - |x| \le \alpha .2^{-52}.|X|$$

$$|X| \le |x| + \alpha \cdot 2^{-52} \cdot |X|$$

According to the collary 3, we have that $ulp(x) \leq 2^{-52}.|x|$ then $ulp(1) = 2^{-52}.$

The notations of the next part which are drawn from [19], we will simplify the calculations of precisions of our algorithm using **unit Roundoff**.

1.2.2 Unit Roundoff

The unit Roundoff has as base 2 and as precisions 53, so we have that:

$$u = \begin{cases} \frac{1}{2} . ulp(1) = \frac{1}{2} . 2^{-52} = 2^{-53} & for RNDN \\ ulp(1) = 2^{-52} & for (RNDU, RNDZ, RNDD) \end{cases}$$

We can do only one calculation and at the end, we can replace u by its values for each rounding mode.

Corollary 4 If $X = \circ(x)$ according to the properties 1, 2, 3 and 4, we have his results:

- $|X x| \le u.|x|$
- $|X x| \le u.|X|$

Corollary 5 If $X = o(x) \Rightarrow |X| \le |x|(1+\epsilon)$ with $|\epsilon| \le u$.

Corollary 6 If
$$X = o(x) \Rightarrow |X| \le |x| + \epsilon$$
 with $|\epsilon| \le u.|X|$.

To proof the exactness of our operation, we use **Sterbenz**'s lemma. The next subsection is taken from [19].

1.2.3 Sterbenz's lemma

Lemma 1 (Sterbenz) In a $radix^1$ - β floating-point system with subnormal numbers available, if x and y are finite floating-point numbers such that $\frac{y}{2} \le x \le 2.y$, then x - y is exactly representable.

 $\beta = 2$ for our research.

Lemma of **Sterbenz** implies for the 4 rounding modes so the result is exact.

PROOF According to the technique of proof of Sterbeinz's lemma from [19]. We suppose that $x \ge 0$, $y \ge 0$ and $y \le x \le 2.y$.

Let $x = M_x \cdot \beta^{e_x - p + 1}$ and $y = M_y \cdot \beta^{e_y - p + 1}$ with their **exponents** (e_x, e_y) and their **mantissa** (M_x, M_y) . We have:

$$\begin{cases} e_{min} \le e_x \le e_{max} \\ e_{min} \le e_y \le e_{max} \\ 0 \le M_x \le \beta^p - 1 \\ 0 < M_y < \beta^p - 1 \end{cases}$$

Firstly, we assume that $e_y \leq e_x$, we define $\lambda = e_x - e_y$. So we have :

$$x - y = (M_x . \beta^{\lambda} - M_y) . \beta^{e_y - p + 1}$$

We define $M = M_x \cdot \beta^{\lambda} - M_y$.

According to the conditions at the beginning of the proof, we have:

- $x \ge y \Rightarrow x y \ge 0 \Rightarrow M.\beta^{e_y p + 1} \ge 0$ as $\beta^{e_y p + 1} > 0 \Rightarrow M \ge 0$.
- $x \le 2.y \Rightarrow x y \le y \Rightarrow M.\beta^{e_y p + 1} \le M_y.\beta^{e_y p + 1} \Rightarrow M \le M_y \le \beta^p 1$

Secondly, we suppose that $M_y < M_x$ and that $e_y > e_x$, we define an other $\lambda = e_y - e_x$.

$$x - y = (M_x - M_y \cdot \beta^{\lambda}) \cdot \beta^{e_x - p + 1}$$

We define $M = M_x - M_y \cdot \beta^{\lambda}$.

According to the conditions at the beginning of the proof, we have:

•
$$x \ge y \Rightarrow x - y \ge 0 \Rightarrow M.\beta^{e_x - p + 1} \ge 0$$
 as $\beta^{e_x - p + 1} > 0 \Rightarrow M \ge 0$.

¹base

• $x \le 2.y \Rightarrow x-y \le y \Rightarrow M.\beta^{e_x-p+1} \le M_y.\beta^{e_y-p+1} \Rightarrow M \le M_y \le \beta^p-1$

We have for our 2 cases, the same result that is $x - y = M.\beta^{e-p+1}$ with $e_{min} \le e_x \le e_{max}$ and $M \le \beta^p - 1$.

We have proved that x - y is a floating point number so the calculation is exact.

1.3 Notations

• Function names are composed of 3 letters followed by 3 numbers. If we have Add then it's an addition and if we have Mul then it's a multiplication.

The first 2 numbers of the 3 represent the argument and the last number represent the output of this function. The number 1 is a **Double** number, 2 is a **Double-Double** and 3 is a **Triple-Double**.

- RNDN, RNDU, RNDD and RNDZ will be rounded as explained before.(table 1.1 and table 1.2)
- \circ () will represent all roundings.
- $cr \log_{fast-path}$, $cr \log_{accurate-path}$ and $cr \log_{advanced-accurate-path}$ for the sake of writing, we will change them to $cr \log_{fast}$, $cr \log_{accurate}$ and $cr \log_{advanced}$.
- Double-Double is a pair of Double numbers, Triple-Double is a triple of Double numbers.

In our research, we use as base $\beta=2,$ precision p=53 , bias=1023 , $E_{min}=-1022$ and as $E_{max}=1023.$

The next chapter will be devited to calculate $cr \log$.

Chapter 2

Steps to calculate crlog

Before starting to implement $cr \log$, we will see what are its steps. According to [13], first, we filter the special cases. Then, we reduce the argument range. After, we use the polynomial approximation and evaluation.

2.1 Special cases

In the special cases, we have five possibilities:

- input is NaN return NaN
- Input is negatif return NaN
- Input is +0 or -0 return $-\infty$
- Input is $+\infty$ return $+\infty$
- Input is a **subnormal** number. We transform it into **normal** number then we calculate as if it is normal.

2.2 Argument reduction

We know that $\log(x.y) = \log(x) + \log(y)$. We have in $Input = 2^E \times m$ with E is the **exposant** and m is the **mantissa**.

So we have $\log(Input) = E \times \log(2) + \log(m)$.

First, we calculate log(m) with **Tang**'s algorithm and **Gal**'s algorithm.

2.2.1Tang's Algorithm

With $x = m * 2^e$ and so that $\log(x) = \log(m) + e * \log(2)$

First, we calculate $\log(m)$ with **Tang**'s Algorithm ([13]).

We have two possibilities: either we make a single reduction or two reductions with this algorithm.

```
# The random() function only gives results that are between 0 and 1.
# We add 1 so that 1 <= m
m = 1 + random(); m
1.672168120969894
```

Figure 2.1: m taken in random to test

2.2.1.1Tang's algorithm with a single reduction

We have : $1 \le m < 2$ and we want to reduce $\log(m)$ using **Tang**'s algorithm. First, we use k significant bits after the initial 1. Then, we take i which represents the integer of the k bits:

So we have: $0 \le i < 2^k$ et $1 + \frac{i}{2^k} \le m < 1 + \frac{i+1}{2^k}$ We look for α_i such as $1 \le m * \alpha_i < 1 + \epsilon_i$ with ϵ_i as small as possible. We can write $m = \frac{m*\alpha_i}{\alpha_i}$ and therefore have $\log(m) = \log(m*\alpha_i) - \log(\alpha_i)$. After calculation, we have $\alpha_i = \frac{2^k}{2^k + i}$ and $\epsilon_i = \frac{1}{2^k + i}$ as $0 \le i < 2^k$ then the value $max(\epsilon_i) = \frac{1}{2^k}$.

If we take $m' = m * \alpha_i$, so we have $1 \leq m' < 1 + \frac{1}{2^k}$. At the end of this algorithm, m' has it's k first zero bits after the initial 1.

The value to be calculated has been reduced to $\log(m) = \log(m')$ $\log(\alpha_i)$ with $-\log(\alpha_i)$ which is already calculated and memorized in a table.

Approach thanks to the Tang's algorithm with a single reduction

```
print("representation of m in binary: ",RR(m).str(2))
binary = RR(m).str(2)[2:9] # k = 7, we take the 7 bits after initial 1 print("the k bits taken: ",binary)
i = ZZ(binary,2) # i is integer of k bits
print("i is integer of k bits: ", i)
# We calculate with the formula alpha_i = (2^k)/(2^k + i)
# We catcutate with the formula atpina_1 = (2 k)/(2 k + 1)
alpha_i = RR(2^k/(2^k + i))
print("alpha_i = ",alpha_i)
# mprime = m*alpha_i and we know that 1 <= mprime < 1+ epsillon_i
# with epsillon_i = 1/(2^k + i)</pre>
epsillon_i = RR(1/(2^k + i))
print("the maximum d'alpha_i not reached : ", epsillon_i)
mprime = m*alpha_i # it can be seen that 1 <= mprime < 1 + epsillon_i
print("mprime = ",mprime)
print("we see the first k bits are set to zero")
print(RR(mprime).str(2))
the k bits taken: 1010110
i is integer of k bits: 86
alpha_i = 0.598130841121495
the maximum d'alpha_i not reached : 0.00467289719626168
mprime = 1.00017532469227
we see the first k bits are set to zero
```

Figure 2.2: The first reduction with k = 7

2.2.1.2 Tang's algorithm with two reductions

After the first reduction explained above, another reduction is made again with the same value of k for the following bits of the first reduction. Then, we take j which represents the integer of the k bits. So we have : $0 \le j < 2^k$ and $1 + \frac{j}{2^{2k}} \le m' < 1 + \frac{j+1}{2^{2k}}$.

```
Now, we are looking for \beta_j such as 1 \leq m' * \beta_j < 1 + \epsilon'_j with \epsilon'_j as small as possible. We have that \log(m') = \log(m' * \beta_j) - \log(\beta_j). After calculation : \beta_j = \frac{2^{2k}}{2^{2k}+j} and \epsilon'_j = \frac{1}{2^{2k}+j}. As we have m'' = m' * \beta_j then 1 \leq m'' < 1 + \frac{1}{2^{2k}}. So, we have \log(m) = \log(m'') - \log(\beta_j) - \log(\alpha_i). (-\log(\alpha_i)) and (-\log(\beta_j)) are already calculated and are put into memory in a table.
```

Approach thanks to the Tang's algorithm with a second reduction.

```
print("representation of mprime in binary: ",RR(mprime).str(2))
k = 7
binary = RR(mprime).str(2)[9:16] # k = 7, we take the 7 bits after initial 1
print("the k bits taken : ",binary)
j = ZZ(binary,2) # j is integer of k bits
print("j is integer of k bits: ", j)
# We calculate with the formula beta_j = (2^2k)/(2^2k + j)
beta_j = RR(2^(2*k)/(2^(2*k) + j))
print("beta_j = ",beta_j)
# msecond = mprime*beta_j and we know that 1 <= msecond < 1+ epsillonprime_j
# avec epsillonprime_j = 1/(2^(2k) + j)
epsillonprime_j = RR(1/(2^(2*k) + j))
print("the maximum d'alpha_i not reached : ", epsillonprime_j)
msecond = mprime*beta_j; msecond # on voit bien 1 <= msecond < 1 + epsillonprime</pre>
 msecond = mprime*beta_j;msecond # on voit bien 1 <= msecond < 1 + epsillonprime_j
print("mprime = ",msecond)
print("we see the first 2k bits are set to zero")
 print(RR(msecond).str(2))
 the k bits taken: 0000010
j is integer of k bits: 2
beta_j = 0.999877944586842
the maximum d'alpha_i not reached : 0.0000610277065787868
mprime = 1.00005324787979
 we see the first 2k bits are set to zero
```

Figure 2.3: The second reduction with k=7

2.2.2Gal's method with double precision

Gal's method is taken from [8].

We have $x = y * 2^n$ with n integer and $0.75 \le y < 1.5$

So $\log(x) = \log(y) + n * \log(2)$.

We calculate $\log(y)$ using a table of triplets $(X_i, \log(X_i), \frac{1}{X_i})$ with $0 \le i \le 1$

 $X_i = 0.75 + \frac{i}{256} + \frac{1}{512} + E_i$ (with E_i a very small number). $\log(y) = \log(\frac{X_i * y}{X_i}) = \log(X_i) + \log(1 + \frac{y - X_i}{X_i}) = \log(X_i) + \log(1 + z)$ If X_i is choosen close to y, we have $\frac{-1}{384} < z < \frac{1}{384}$ and if more y is close to 1 so $\frac{-1}{512} < z < \frac{1}{512}$.

We use an approximation polynomial p(z) of degree 6 (for double precision) to approach the function $\log(1+z)$.

If x is close to 1 so we use an approximation polynomial of $\log(x)$ without the table with a relative error of 2^{-72} which is negligible.

The table of the triplets $(X_i, \log(X_i), \frac{1}{X_i})$ contains 576 elements.

$$F_i = \log(X_i)$$
 and $G_i = \frac{1}{X_i}$ with $0 \le i \le 192$.

The numbers X_i are choosen so that F_i and G_i are 56 bits of mantissa. And they have a relative precision of 2^{-65} .

We search to be close to $0.75 + \frac{i}{256} + \frac{1}{512}$ with X_i such that bits 57 to 67 of the mantissa of $\log(X_i)$ and of $\frac{1}{X_i}$ which will be reset either all to 0 or

all to 1. That's why small numbers E_i were introduced. F_i and G_i are the numbers with double precision obtained by a calculation of extended precisions and a symmetric rounding.

2.2.3 method with Tang algorithm and Gal method

Our routine involves the two algorithms seen previously.

We start using **Tang**'s algorithm.

at the beginning as explained in 2.2.1, we only have $x = m*2^e$ and therefore we have $\log(x) = \log(m) + e * \log(2)$.

As mentioned in 2.2.1, we recover α_i with the calculation:

 $\alpha_i = \frac{2^k}{2^k + i}$ except for the following we will take α_i' a double close to α_i such that the $\log(\alpha_i')$ is the closest to a double.

This ϵ_i will be used to search for the $\log(\alpha_i')$ such that the bits from 54^{th} to 71^{st} are identical. (To have an approximation approximately 2^{-71}). (Figure 2.4)

method with Tang's algorithm and the Gal method

```
k=8 # we take k = 8,
m=1+random() # we take 1<m<2
M = R200(m).exact_rational()
binary = RR(m).str(2)[2:10] #we recover the 8 bits after the initial 1.
i = int(binary,2) # i is the 8-bit integer
alpha_i_m = R200(table_alpha_modified[i],16).exact_rational() #table calculated for all i of alpha_i_m
                                            # such that log(alpha_i_m) has a precision of 71 bits
log_alpha_i_m= R200(table_log_alpha_modified[i],16).exact_rational()
A=R200(log(M*alpha_i_m)+log_alpha_i_m) #
B = R200(log(R200(M)))
print("m =",m)
print("i =",i)
print ("alpha_iprime = ",alpha_i_m )
print("log_alpha_i_prime :",log_alpha_i_m)
print("sol =" ,A)
print("sol =" ,A
print("log(m)=",B)
C = R200((A.exact_rational()-B.exact_rational())/B.exact_rational())
print("the precision for the modified method is %f bits" %(R200(-log(abs(C))/log(2.0))))
m = 1.284464826386863
alpha_iprime = 7030009188575477/9007199254740992
log_alpha_i_prime : 8929238770791491/36028797018963968
       = 0.25034215407638279200093780192379456668445060635973530859479
log(m) = 0.25034215407638279200088487247857236200406769824813103788058
the precision for the modified method is 72.002249 bits
```

Figure 2.4: The revised method of **Tang** and **Gal**

We are going to verify if we obtain the same result with the **Tang** algorithm and our modified method (Figure 2.5).

Verification of our algorithm with that of Tang

```
k=8 # we take k=8
  m=1+random() # we take 1<m<2
 M = R200(m).exact_rational()
 binary = RR(m) .str(2) [2:10]
 i = int (binary,2)# we recover the 8 bits after the initial 1.
(h1,l1) = table_alpha_i[i] # h1 is the main value and l1 is the error value
 (H1,L1)= (R200(h1,16).exact_rational(),R200(l1,16).exact_rational())
 (h2, l2)= table_log_alpha_i[i]# h2 is the main value and 12 is the error value
(H2,L2)= (R200(h2,16).exact_rational(),R200(l2,16).exact_rational())
 alpha_i_m = R200(table_alpha_modified[i],16) #table calculated for all i of alpha_i_m
# such that log(alpha_i_m) has a precision of 71 bits
 log_alpha_i_m= R200(table_log_alpha_modified[i],16)
sol_Tang = R200(log(M*(H1+L1))+(H2+L2))
sol_rema=R200(log(M*alpha_i_m)+log_alpha_i_m)
 A = R200(log(R200(M)))
print("log(m) =
C = R200((sol_Tang.exact_rational()-A.exact_rational())/A.exact_rational())
D = R200((sol_rema.exact_rational()-A.exact_rational())/A.exact_rational())
print("the precision for Tang's method is %f bits" %(R200(-log(abs(C))/log(2.0))))
print("the precision for the reworked method is %f bits" %(R200(-log(abs(D))/log(2.0))))
 m = 1.486258461986336
 i = 124

'alpha,i.m = ('0x1.58ed2308158edp-1', '0x1.947941c2116fbp-2')
 alpha_1_m : 0.67368421158545022109365163487382233142852783203125000000000
 solution of Tang's method = 0.39626186252142292454373873382587912813786396912842268809010
 solution of the reworked method :
                                                                  0.39626186252142292454384333370528885498555363956773361536439
                                              0.39626186252142292454373873382587662034589290447866241791757
 the precision for Tang's method is 106.961735 bits
 the precision for the reworked method is 71.682063 bits
```

Figure 2.5: Verification of the revised method and **Tang**'s method

2.3 polynomial approximation and evaluation

Before calculating the approximation function with the **Sollya** tool, we will look for the interval for which the polynomial will be effective for $m*\alpha_i'$ with $0 \le i < 256$ and $1 + \frac{i}{256} \le m < 1 + \frac{i+1}{256}$. The calculations are experimented on **sage**, see the diagram 2.6:

Calculation of the interval for the approximation polynomial

```
Linf = [] # the lower bounds for each m*alpha_i_m
Lsup = [] # the upper bounds for each m*alpha_i_m
for i in range(256):
    Linf.append(R200((1+i/256)*R200(table_alpha_modified[i],16).exact_rational()))
    Lsup.append(R200((1+(1+i)/256)*R200(table_alpha_modified[i],16).exact_rational()))
print("the upper bounds is : ", max(Lsup).exact_rational())
print("the lower bounds is : ", min(Linf).exact_rational())

the upper bounds is : 257/256
the lower bounds is : 1
```

Figure 2.6: Interval calculation for the approximation polynomial

We find as result $1 < m.\alpha_i' < \frac{257}{256}$, exactly the same bounds as we calculated with the **Tang** method.

At the end of the calculation of the modified algorithm from Tang and Gal's algorithm, we have that $\log(input) = \log(m.\alpha_i) - \log(\alpha_i) + E.\log(2)$. We transform $\log(m.\alpha_i)$ into $\log(1+t)$ to use the approximation function; then we have $t = m.\alpha_i - 1$.

We search the polynomial approximation for $\log(1+t)$ thanks to the **Taylor** formula.

Let P(t) the polynomial approximation for $\log(1+t)$, we have $P(t) \approx t - \frac{t^2}{2}$..., in case we have a constant term, we can reduce it to 0.

This calculation will be done thanks to the **Sollya** tool with its **fpminimax** function with the calculated interval.

Then we evaluate this approximation function with the argument $t = m.\alpha_i - 1$.

Now, we have the operation $\log(input) = P(m.\alpha_i - 1) - \log(\alpha_i)$ (Figure 2.7).

Our method with the Tang algorithm and the Gal method with the P function

```
k=8 # we take k = 8
M = R200(m).exact_rational()
binary = RR(m).str(2)[2:10] # we recover the 8 bits after the initial 1.
i = int(binary,2) # i is the 8-bit integer
alpha_i_m = R200(table_alpha_modified[i],16).exact_rational() #table calculated for all i of alpha_i_m
                                          # such that log(alpha_i_m) has a precision of 71 bits
log_alpha_i_m= R200(table_log_alpha_modified[i],16).exact_rational()
A=R200(P(M*alpha_i_m-1)+log_alpha_i_m) #
B = R200(log(R200(M)))
print("m =",m)
print("i =",i)
print("alpha_iprime = ",alpha_i_m )
print("log_alpha_i_prime :",log_alpha_i_m)
print("sol =" ,A)
print("sol =" ,A
print("log(m)=",B)
C = R200((A.exact_rational()-B.exact_rational())/B.exact_rational())
print("the precision for the reworked method is %f bits" %(R200(-log(abs(C))/log(2.0))))
m = 1.019019152905743
alpha_iprime = 8868626965695373/9007199254740992
log_alpha_i_prime : 8937554567599383/576460752303423488
      = 0.018840549849761655851050549366714043792850403560626948245975
log(m)= 0.018840549849761655851046655024229337473701007314919333872294
the precision for the reworked method is 72.034879 bits
```

Figure 2.7: The revised method of **Tang** and **Gal** with P(t)

Before talking about the implementation of cr_{log} , we will first see the used algorithms.

The algorithms of the chapter 4 and chapter 5 are taken from [4] and [12].

Chapter 3

Operators on Double-Double numbers

3.1 Addition Operators

3.1.1 Add112

See algorithm 1

Lemma 1 (Add112) Let a and b floating point numbers, with $|a| \ge |b|$, s and t result of Add112(a,b) for the 4 modes of rounding, considering that there is no **overflow** so:

- (1) s + t is exactly equal to a + b.
- (2) $|t| \le 2^{-53}|s|$ for RNDN and $|t| \le 2^{-52}|s|$ for direct rounding modes.

Proof (1):

According to the calculation technique of **Fast2Sum algorithm** proof's ([19]), this proof shown that for the closest rounding mode (RNDN).

We suppose that $|a| \ge |b|$ and there is no **overflow**.

We suppose a > 0 and b > 0 (respectively a < 0 and b < 0).

We take $b = b_h + b_\ell$ with b_h a multiple of ulp(a), $|b_\ell| < ulp(a)$ and b_ℓ is of the same sign that b (and b_h).

Either we have $0 \le |b_{\ell}| < 2^{-53}.|a|$ or $2^{-53}.|a| \le |b_{\ell}| \le 2^{-52}.|a|$. If $0 \le |b_{\ell}| < 2^{-53}.|a|$ for the modes of Rounding RNDZ, RNDN and RNDD (respectively RNDZ, RNDN and RNDU) then b_{ℓ} will be ignored otherwise RNDU(respectively RNDD) b_{ℓ} is not ignore, so $\circ(b)$ $b_h + ulp(a);$

If $2^{-53} \cdot |a| \leq |b_{\ell}| \leq 2^{-52} \cdot |a|$ for the modes of Rounding RNDZ and RNDD (respectively RNDZ and RNDU) then b_{ℓ} will be also ignored for RNDN (closest round) and RNDU (respectively RNDN and RNDD).

The first case when b_{ℓ} is ignored so we have $\circ(b) = b_h$. \Rightarrow :

- $s = \circ(a+b) \Rightarrow s = a+b_h$
- $z = \circ(s-a) \Rightarrow z = a + b_h a \Rightarrow z = b_h$
- $t = \circ(b-z) \Rightarrow t = b-b_h$ we have that $b = b_h + b_\ell \Rightarrow t = b_\ell$

So we have that $s + t = a + b_h + b_\ell \Rightarrow s + t$ is exactly equal to a + b.

The second case when b_{ℓ} is not ignored so we have $\circ(b) = b_h + ulp(a)$. \Rightarrow :

- $s = \circ(a+b) \Rightarrow s = a+b_h+ulp(a)$
- $z = \circ(s-a) \Rightarrow z = a + b_h + ulp(a) a \Rightarrow z = b_h + ulp(a)$
- $t = \circ(b-z) \Rightarrow t = b (b_h + ulp(a))$ we have that $b = b_h + b_\ell \Rightarrow t = b_\ell ulp(a)$

So we have that $s + t = a + b_h + ulp(a) + b_\ell - ulp(a) \Rightarrow s + t$ is exactly equal to a + b.

We suppose a>0 and b>0 (respectively a<0 and b<0). We suppose that |b|< ulp(a). Either we have $0<|b|<2^{-53}.|a|$ or $2^{-53}|a|\leq |b|\leq 2^{-52}.|a|$.

If $0 < |b| < 2^{-53}.|a|$ for the modes of Rounding RNDZ, RNDN and RNDD (respectively RNDZ, RNDN and RNDU) then b will be ignored otherwise RNDU(respectively RNDD) b is not ignore, so we have $\circ(b) = ulp(a)$;

If $2^{-53}|a| \leq |b| \leq 2^{-52}.|a|$ for the modes of Rounding RNDZ and RNDD (respectively RNDZ and RNDU) then b_{ℓ} will be also ignored for RNDN (closest round) and RNDU (respectively RNDN and RNDD).

The first case when b is ignored so we have :

- $s = \circ(a+b) \Rightarrow s = a$
- $z = \circ(s a) \Rightarrow z = a a \Rightarrow z = 0$
- $t = \circ(b-z) \Rightarrow t = b 0 \Rightarrow t = b$

So we have that $s + t = a + b \Rightarrow s + t$ is exactly equal to a + b.

The second case when b is not ignored so we have $\circ(b) = ulp(a)$.

- $s = \circ(a+b) \Rightarrow s = a + ulp(a)$
- $z = \circ(s a) \Rightarrow z = a + ulp(a) a \Rightarrow z = ulp(a)$
- $t = \circ(b-z) \Rightarrow t = b ulp(a)$

So we have that $s+t=a+ulp(a)+b-ulp(a) \Rightarrow s+t$ is exactly equal to a+b.

We suppose that a>0 and b<0 (respectively a<0 and b>0): If $|b|\geq |\frac{a}{2}|$.

So we have:

- s = o(a+b), we have $\frac{-b}{2} \le a \le -2b$ after the **Sterbenz**'s lemma (lemma 1), we have s is exactly equal to a+b
- $z = \circ(s-a) \Rightarrow z = \circ((a+b)-a) \Rightarrow z = \circ(b) = b \Rightarrow z = b$
- $t = \circ(b-z) \Rightarrow t = \circ(b-b) \Rightarrow t = 0$

So we have, s + t is exactly equal to a + b.

If $|b| < |\frac{a}{2}|$:

- $s = \circ(a+b) \Rightarrow \frac{a}{2} < s \le a \Rightarrow \frac{a}{2} \le s \le 2a$
- $z = o(s a) \Rightarrow \text{ and } \frac{a}{2} \le s \le 2a$ after the **Sterbenz**'s lemma (lemma 1), we have z is exactly equal to s a.
- $t = \circ(b-z) \Rightarrow t = \circ(b-(s-a)) \Rightarrow t = \circ(a+b-s) \Rightarrow t = \circ(a+(b-s))$ or $-\frac{a}{2} < b < 0$ and $\frac{a}{2} \le s \le 2a \Rightarrow \frac{a}{2} 0 < s b < 2a (-\frac{a}{2}) \Rightarrow \frac{a}{2} < s b < \frac{3a}{2} \Rightarrow \frac{a}{2} \le s b \le 2a$ after the **Sterbenz**'s lemma (lemma 1), we have t is exactly equal to a + b s.

We have $s + t = a + b - s + s \Rightarrow s + t = a + b$.

(By symetric, we have the same results for a < 0 and b > 0.)

The 4 modes of rounding for Add112 have an exact equality between s + t and a + b.

(2):

According to Proof(1), We have 6 possibilities for the result of s and t:

- If s = a + b and t = 0 so $|t| \le 2^{-53}|s|$ for RNDN and $|t| \le 2^{-52}|s|$ for the direct rounding modes.
- If s = a and t = b (a and b have the same sign) with $|b| < 2^{-53}$. $|a| \Rightarrow$ $|t| \leq 2^{-53}|s|$ for RNDN a fortiori $|t| \leq 2^{-52}|s|$ for the direct rounding modes.
- If $s = a + b_h$ and $t = b_l$ (a, b_h and b_ℓ are the same sign). We suppose that $a > 0, b_h > 0$ and $b_\ell > 0$. We know that $0 \le b_\ell < 2^{-53}.a \Rightarrow$ $a > 2^{53}.b_{\ell}$ and that b_h is a multiple of $ulp(a) \Rightarrow b_h > 2.b_{\ell}$. So we have $a + b_h > 2^{53}b_\ell + 2b_\ell \Rightarrow a + b_h > (2^{53} + 2).b_\ell$ or $2^{53} + 2 > 2^{53}$ $\Rightarrow a + b_h > 2^{53}b_{\ell} \Rightarrow s > 2^{53}.t \Rightarrow t \leq 2^{-53}.s \Rightarrow |t| \leq 2^{-53}.|s|$ (and similarly if a < 0, $b_h < 0$ and $b_\ell < 0$).
- If s = a + ulp(a) and t = b ulp(a): We know according to proof(1) that $|b| \leq 2^{-53}|a|$. We suppose that a > 0 and b > 0, so we have: $t \le t + ulp(a) \Rightarrow t \le b - ulp(a) + ulp(a) \Rightarrow t \le b \Rightarrow |t| \le 2^{-53}|a|$ We search the upper bound of |a| as a function of |s|: $s = a + ulp(a) \Rightarrow s \ge a \Rightarrow |s| \ge |a|$ so: $|t| \le 2^{-53}|s|$ (and similarly if a < 0 and b < 0).
- If $s = a + b_h + ulp(a)$ and $t = b_\ell ulp(a)$ We know according to proof(1) that $|b_{\ell}| \leq 2^{-53} |b_h|$. We suppose that a > 0 and $b_h > 0$, so we have: $t \le t + ulp(a) \Rightarrow t \le b_{\ell} - ulp(a) + ulp(a) \Rightarrow t \le b_{\ell} \Rightarrow |t| \le 2^{-53}|b_h|$ $\Rightarrow |t| < 2^{-53}|a + b_h|$ We search the upper bound of $|a + b_h|$ as a function of |s|: $s = a + b_h + ulp(a) \Rightarrow s \ge a + b_h \Rightarrow |s| \ge |a + b_h|$ so: $|t| < 2^{-53}|s|$ (and similarly if a < 0 and b < 0).

• If t = a + b - s, so |t| = |a + b - s| = |s - (a + b)| according to collary $6 \Rightarrow$:

$$|s - (a+b)| \le u.|s|$$

with u is the unit Roundoff (see 1.2.2).

For all cases, we have $|t| \le u.|s|$ so for RNDN: $|t| \le 2^{-53}.|s|$ and for the others rounding modes: $|t| \le 2^{-52}.|s|$.

$3.1.2 \quad Add122$

See algorithm 2

Lemma 2 (Add122) Let a a **Double** number and (b_h, b_ℓ) a **Double-Double** number, with $|a| \ge |b_h|$, s and t result of $Add122(a, b_h, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so:

- (1) $s + t = a + b_h + b_\ell + \theta \text{ with } |\theta| \le u.|t|$
- $(2) |t| \leq 2.u.|s|$

 $2.u = 2^{-52}$ for RNDN and $2.u = 2^{-51}$ for the other rounding modes.

PROOF (1) We suppose $|a| \ge |b_h|$ and there is no **overflow**. We have:

- $s + \ell = Add112(a, b_h)$ according to Lemma $1 \Rightarrow s + \ell$ is exactly equal to $a + b_h$.
- $t = o(\ell + b_{\ell}) \Rightarrow t = \ell + b_{\ell} + \theta$ according to the collary 6 we have $|\theta| \le u.|t|$

So, we have $s+t=a+b_h-\ell+l+b_\ell+\theta_1 \Rightarrow s+t=a+b_h+b_\ell+\theta_1 \Rightarrow s+t=a+b+\theta$ with $|\theta| \leq u.|t|$.

- (2) We suppose $|a| \ge |b_h|$ and there is no **overflow**. We have:
 - $t = o(\ell + b_{\ell}) \Rightarrow t = \ell + b_{\ell} + \epsilon$ with $|\epsilon| < u.|t|$

• $s + \ell = Add112(a, b_h)$ according to Lemma $1 \Rightarrow s + \ell = a + b_h$ and $|\ell| \leq u.|s| \Rightarrow$

$$|s + \ell| \ge |s| - |\ell|$$
$$|s + \ell| \ge |s| - u.|s|$$
$$|s + \ell| \ge (1 - u).|s|$$

 \Rightarrow

$$(1-u).|s| \le |a+b_h|$$

$$|s| \le \frac{1}{1-u}.|a+b_h|$$

As
$$|\ell| \le u.|s| \Rightarrow$$

$$|\ell| \le \frac{u}{1-u}.|a+b_h|$$

As $|b_{\ell}| \leq u.|b_h|$ and $|a| \geq |b_h| \Rightarrow |b_{\ell}| \leq u.|a|$ and that $|a| \leq 2.|a+b_h| \Rightarrow$

$$|b_{\ell}| \le 2.u.|a + b_h|$$

$$|t| \le \frac{u}{1 - u}.|a + b_h| + 2.u.|a + b_h| + u.|t|$$

$$(1 - u).|t| \le (\frac{u}{1 - u} + 2.u).|a + b_h|$$

$$|t| \le \frac{1}{1 - u}.(\frac{u}{1 - u} + 2.u).|a + b_h|$$

We seek the upper bound of $|a + b_h|$ in function of |s|:

$$|a+b_h| \le (1+u).|s|$$

 \Rightarrow

$$|t| \le \frac{1}{1-u} \cdot \left(\frac{u}{1-u} + 2 \cdot u\right) \cdot (1+u) \cdot |s|$$
$$|t| \le \frac{(-2 \cdot u^2 + 2 \cdot u) \cdot (u+1)}{(1-u)^2} \cdot |s|$$

As $\frac{(-2.u^2+2.u).(u+1)}{(1-u)^2} \le 2.u$ We have:

$$|t| \le 2.u.|s|$$

Theorem 1 (Relative error algorithm Add122 without occurring of cancellation)

Let a a **double** number and (b_h, b_ℓ) a **double-double** number are the arguments of the function Add122:

If a and b_h are the same sign, so:

$$s + t = (a + b_h + b_\ell)(1 + \epsilon)$$

with $|\epsilon| \le 2.u^2$.

 $2.u^2 = 2^{-105}$ for RNDN and $2.u^2 = 2^{-103}$ for the other rounding modes.

PROOF According to the calculation technique of the theorem 4.2 proof's([12]).

We have $|a| \ge |b_h|$, either a > 0 and $b_h > 0$ or a < 0 and $b_h < 0$. As they're symmetric, we only use a > 0 and $b_h > 0$.

Based on the algorithm (Add122), we have:

 $t = o(l + b_{\ell})$ according to collary $5 \Rightarrow t = (l + b_{\ell}).(1 + \epsilon)$ with $|\epsilon| \leq u$.

$$t = \ell + b_{\ell} + \delta$$

with $\delta = (\ell + b_{\ell}).\epsilon$ We calculate $|\delta|$, so we have:

$$|\delta| = |(\ell + b_{\ell}).\epsilon$$

according to the triangle inequality, we have:

$$|\delta| \le |\ell.\epsilon| + |b_{\ell}.\epsilon|$$

$$|\delta| \le |\epsilon|.(|\ell| + |b_{\ell}|)$$

Based on Lemma 1, we have $|l| \le u.|s| \Rightarrow |s + \ell| \ge (1 - u).|s| \Rightarrow (1 - u).|s| \le |a + b_h| \Rightarrow |s| \le \frac{1}{1 - u}.|a + b_h| \Rightarrow |\ell| \le \frac{u}{1 - u}.|a + b_h|$

So we have:

$$|\delta| \le |\epsilon| \cdot \left(\frac{u}{1-u} \cdot |a+b_h| + |b_\ell|\right)$$

we have $|b_{\ell}| \leq u.|b_h| \leq u.|a| \leq u.|a+b_h| \Rightarrow$

$$|\delta| \le |\epsilon| \cdot \left(\frac{u}{1-u} \cdot |a+b_h| + u \cdot |a+b_h|\right)$$

$$|\delta| \le |\epsilon| \cdot \frac{-u^2 + 2 \cdot u}{1 - u} \cdot |a + b_h|$$

As $\frac{-u^2+2.u}{1-u} \leq 2.u$ and $|\epsilon| \leq u \Rightarrow$

$$|\delta| \le 2.u^2.|a + b_h|$$

$$|b_{\ell}| \leq u.|a|$$

$$|b_{\ell}| \leq u.|a+b_h|$$

Now we search a lower bound for $|a + b_h| + b_\ell$ in function $|a + b_h|$

$$|a + b_h + b_\ell| \le |a + b_h| + |b_\ell|$$

So we have:

$$|a + b_h + b_\ell| \ge (1 - u)|a + b_h|$$

and we knows that $|\delta| \leq |a + b_h| \times 2.u^2 \Rightarrow$

$$|\delta| \le |a + b_h + b_\ell| (\frac{1}{1 - u}).2.u^2$$

After the calculations, we have:

$$|\delta| \le |a + b_h + b_\ell|.|\epsilon|$$

with $|\epsilon| \leq \frac{2.u^2}{1-u} \leq 2.u^2$

Theorem 2 (Relative error algorithm Add122 with a bounded cancellation)

Let a a **Double** number and (b_h, b_ℓ) a **double-double** number. We have for the algorithm Add122 with a and (b_h, b_ℓ) it's arguments. If a and b_h are different sign and we suppose $|b_h| \leq 2^{-\mu} |a|$ with $\mu \geq 1$. So:

$$s + t = (a + (b_h + b_\ell))(1 + \epsilon)$$

with

$$|\epsilon| \le 2.u^2 \cdot \frac{2^{-53-\mu}}{1-2^{-\mu}-u} \le 2.(2.u^2) = 4.u^2$$

PROOF According to the calculation technique of the theorem 4.3 proof's([12]). We suppose $|b_h| \leq 2^{-\mu}|a|$ with $\mu \geq 1$. First we look for the upper bound, using the results of the proof of Theorem 1. We have:

$$|b_{\ell}| \le u.|b_h| \le u.2^{-\mu}|a|$$

then we search the lower bound for $|a + b_h|$ in function of |a|:

$$|a+b_h| > (1-2^{-\mu}).|a|$$

So we have:

$$|b_{\ell}| \le \frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu}} |a + b_h|$$

So:

$$|\delta| \le |a_h + b_h| \cdot 2 \cdot u^2 \cdot \frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu}}$$

Now, we search the lower bound for $|a + b_h + b_\ell|$ depending on $|a + b_h|$. We have :

$$|b_{\ell}| \le \frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu}} |a + b_h|$$
$$|a + b_h + b_{\ell}| \ge |a + b_h| \frac{1 - 2^{-\mu} - u}{1 - 2^{-\mu}}$$

So we have for $|\delta|$:

$$|\delta| \le |a + b_h + b_\ell| \cdot \frac{1 - 2^{-\mu}}{1 - 2^{-\mu} - u} \cdot 2 \cdot u^2 \cdot \frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu}}$$
$$|\delta| \le |a + b_h + b_\ell| \cdot 2 \cdot u^2 \cdot \frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu} - u}$$

So:

$$|\epsilon| \le 2.u^2 \cdot \frac{2^{-53-\mu}}{1 - 2^{-\mu} - u}$$

As $\mu \geq 1$, we want the upper bound of $\frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu} - u} \Rightarrow$

$$\frac{2^{-54}}{1/2 - 2^{-52}} \le 2$$

because $2^{-53} \le u \le 2^{-52}$ So :

$$|\epsilon| \le 2.2.u^2 = 4.u^2$$

3.1.3 Add222

See algorithm 3

Lemma 3 (Add222) Let (a_h, a_ℓ) and (b_h, b_ℓ) **Double-Double** numbers, with $|a_h| \ge |b_h|$, s and t result of $Add222(a_h, a_\ell, b_h, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so:

 $6.u = 2^{-50.4}$ for RNDN and $6.u = 2^{-49.4}$ for the other rounding modes.

PROOF According to $Add222 \Rightarrow : t = o(m + b_{\ell})$ (corrollary 5) $\Rightarrow t \leq (m + b_{\ell})(1 + \epsilon_1)$ with $|\epsilon_1| \leq u$.

$$|t| \le (|m| + |b_{\ell}|).|1 + \epsilon_1|$$

We have: $m = o(\ell + a_{\ell}) \Rightarrow m = (\ell + a_{\ell}).(1 + \epsilon_2)$ with $|\epsilon_2| < u \Rightarrow$.

$$|t| \le (|(\ell + a_{\ell})(1 + \epsilon_2)| + |b_{\ell}|).|1 + \epsilon_1|$$

After calculate, we have:

$$|t| \le |\ell|.|1 + \epsilon_1 + \epsilon_2 + \epsilon_1.\epsilon_2| + |a_{\ell}|.|1 + \epsilon_1 + \epsilon_2 + \epsilon_1.\epsilon_2| + |b_{\ell}|.|1 + \epsilon_1|$$

As $|\epsilon_1| \le u$ and $|\epsilon_2| \le u \Rightarrow$

$$|t| \le |\ell| \cdot (1 + 2 \cdot u + u^2) + |a_{\ell}| \cdot (1 + 2 \cdot u + u^2) + |b_{\ell}| \cdot (1 + u)$$

As:

$$|a_{\ell}| \le u.|a_h| \le 2.u.|a_h + b_h|$$

 $|b_{\ell}| \le u.|b_h| \le u.|a_h| \le 2.u.|a_h + b_h|$

and that $a_h + b_h = s + \ell \Rightarrow$

$$|t| \le |\ell| \cdot (1 + 2 \cdot u + u^2) + 2 \cdot u \cdot (1 + 2 \cdot u + u^2) \cdot |s + \ell| + 2 \cdot u \cdot (1 + u) \cdot |s + \ell|$$

We know that $|\ell| \leq u.|s|$, after calculate, we have:

$$|t| \le (u^4 + 5.u^3 + 7.u^2 + 3.u).|s|$$

As $(2.u^4 + 9.u^3 + 12.u^2 + 5.u) \le 6.u$

$$|t| \leq 6.u.|s|$$

Theorem 3 (Relative error algorithm Add222 without occurring of cancellation) Let (a_h, a_ℓ) and (b_h, b_ℓ) the **double-double**. We have for the algorithm Add222 with (a_h, a_ℓ) and (b_h, b_ℓ) it's arguments. If a_h and b_h have the same sign, so:

$$s + t = ((a_h + a_\ell) + (b_h + b_\ell))(1 + \epsilon)$$

with $|\epsilon| < 3.u^2$

 $3.u^2 = 2^{-104.4}$ for RNDN and $3.u^2 = 2^{-102.4}$ for the other rounding modes.

PROOF According to the calculation technique of the theorem 4.2 proof's([12]). We have $|a_h| \ge |b_h|$, either $a_h > 0$ and $b_h > 0$ or $a_h < 0$ and $b_h < 0$. As they're symmetric, we only use $a_h > 0$ and $b_h > 0$.

Based on the algorithm (Add222), we have:

 $t = o(m + b_{\ell})$ according to corollary $5 \Rightarrow t = (m + b_{\ell})(1 + \epsilon_1)$ with $|\epsilon_1| \leq u$.

$$t = (\circ(\ell + a_{\ell}) + b_{\ell})(1 + \epsilon_1)$$

$$t = ((\ell + a_{\ell})(1 + \epsilon_2) + b_{\ell})(1 + \epsilon_1)$$

with $|\epsilon_2| \leq u$.

$$t = \ell + a_{\ell} + b_{\ell} + \delta$$

with $\delta = (\ell + a_{\ell} + b_{\ell})\epsilon_1 + (\ell + a_{\ell})\epsilon_1\epsilon_2$.

We calculate $|\delta|$, so we have:

$$|\delta| = |(\ell + a_{\ell} + b_{\ell})\epsilon_1 + (\ell + a_{\ell})\epsilon_1\epsilon_2|$$

according to the triangle inequality, we have:

$$|\delta| \le |(\ell + a_{\ell} + b_{\ell})\epsilon_1| + |(\ell + a_{\ell})\epsilon_1\epsilon_2|$$

$$|\delta| \le |\ell\epsilon_1| + |\ell\epsilon_1\epsilon_2| + |(a_{\ell} + b_{\ell})\epsilon_1| + |a_{\ell}\epsilon_1\epsilon_2|$$

$$|\delta| \le |\ell|(|\epsilon_1| + |\epsilon_1\epsilon_2|) + |(a_{\ell} + b_{\ell})\epsilon_1| + |a_{\ell}\epsilon_1\epsilon_2|$$

$$|\delta| \le |\ell| \cdot (|\epsilon_1| + |\epsilon_1 \epsilon_2|) + |(a_\ell + b_\ell)| \cdot |\epsilon_1| + |a_\ell| \cdot |\epsilon_1 \epsilon_2|$$

$$|\delta| \le |\ell| \cdot (|\epsilon_1| + |\epsilon_1 \epsilon_2|) + |(a_\ell + b_\ell)| \cdot |\epsilon_1| + |(a_\ell + b_\ell)| \cdot |\epsilon_1 \epsilon_2|$$

$$|\delta| \le (|\ell| + |(a_\ell + b_\ell)|) \cdot |\epsilon_1 + \epsilon_1 \epsilon_2|$$

Based on Lemma 1, we have $|l| \le u.|s|$ and as $|s| \le |a_h + b_h| \Rightarrow |l| \le u.|a_h + b_h|$.

$$|\delta| \le (u.|a_h + b_h| + |a_\ell| + |b_\ell|).|\epsilon_1 + \epsilon_1 \epsilon_2|$$
we have $|a_\ell| \le u.|a_h| \le u.|a_h + b_h|$ and $|b_\ell \le u.|b_h| \le u.|a_h + b_h| \Rightarrow$

$$|\delta| \le (u.|a_h + b_h| + u.|a_h + b_h| + u.|a_h + b_h|).|\epsilon_1 + \epsilon_1 \epsilon_2|$$

$$|\delta| \le |a_h + b_h| \times 3.u.|\epsilon_1 + \epsilon_1 \epsilon_2|$$

$$|\delta| \le |a_h + b_h| \times 3.u.(u + u^2)$$

$$|\delta| \le |a_h + b_h| \times 3.(u^2 + u^3)$$

We seek the upper bound of $|a_h + b_h|$ in function of $|a_h + b_h + a_\ell + b_\ell|$:

$$|a_{\ell} + b_{\ell}| \le |a_{\ell}| + |b_{\ell}|$$

$$|a_{\ell} + b_{\ell}| \le u.|a_{h}| + u.|b_{h}|$$

$$|a_{\ell} + b_{\ell}| \le 2.u.|a_{h}|$$

$$|a_{\ell} + b_{\ell}| \le 2.u.|a_{h} + b_{h}|$$

 \Rightarrow

$$|a_h + b_h + a_\ell + b_\ell| \ge (1 - 2.u)|a_h + b_h|$$

and we know that $|\delta| \leq |a_h + b_h| \times 3.(u^2 + u^3) \Rightarrow$

$$|\delta| \le |a_h + a_l + b_h + b_\ell| (\frac{1}{1 - 2.u}) 3.(u^2 + u^3)$$

$$|\delta| \le |a_h + a_l + b_h + b_\ell| (\frac{3.u^2 + 3.u^3}{1 - 2.u})$$

but $\frac{3.u^2+3.u^3}{1-2.u} \le 3.u^2$ we have :

$$|\epsilon| \le 3.u^2$$

Theorem 4 (Relative error algorithm Add222 with a bounded cancellation) Let (a_h, a_ℓ) and (b_h, b_ℓ) are the **double-double** number. We have for the algorithm Add222 with (a_h, a_ℓ) and (b_h, b_ℓ) it's arguments.

If a_h and b_h are different sign and we suppose $|b_h| \le 2^{-\mu} |a_h|$ with $\mu \ge 1$. So:

$$s+t=((a_h+a_\ell)+(b_h+b_\ell))(1+\epsilon)$$

with

$$|\epsilon| \le 3.u^2 \cdot \frac{1 - 2^{-\mu - 1}}{1 - 2^{-\mu} - 2.u} \le 2.3.u^2 = 6.u^2$$

PROOF According to the calculation technique of the theorem 4.3 proof's([12]). We suppose $|b_h| \leq 2^{-\mu} |a_h|$ with $\mu \geq 1$. First we look for the upper bound, using the results of the proof of Theorem 4.1.1. We have:

$$|b_{\ell}| \le u.|b_h| \le u.2^{-\mu}|a_h|$$

then we search the lower bound:

$$|a_h + b_h| \ge (1 - 2^{-\mu}).|a_h|$$

3.1. ADDITION OPERATORS

39

As $|a_{\ell}| \le u.|a_h|$ and $|b_{\ell}| \le u.|b_h| \le u.2^{-\mu}.|a_h|$.

$$|a_{\ell}| \le \frac{u}{1 - 2^{-\mu}} |a_h + b_h|$$

and

$$|b_{\ell}| \le \frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu}} |a_h + b_h|$$

 \Rightarrow

$$u + u \cdot 2^{-\mu} \le 1 - 2^{-\mu - 1}$$

because of $\mu \geq 1$

So:

$$|\delta| \le |a_h + b_h| \cdot 3 \cdot u^2 \cdot \frac{1 - 2^{-\mu - 1}}{1 - 2^{-\mu}}$$

Now, we search the lower bound for $|a_h + a_\ell + b_h + b_\ell|$ depending on $|a_h + b_h|$. We have:

$$|a_{\ell} + b_{\ell}| \le |a_{\ell}| + |b_{\ell}|$$

$$|a_{\ell} + b_{\ell}| \le 2.u.|a_{h}|$$

$$|a_{\ell} + b_{\ell}| \le \frac{2.u}{1 - 2^{-\mu}}|a_{h} + b_{h}|$$

$$|a_{h} + a_{\ell} + b_{h} + b_{\ell}| \ge |a_{h} + b_{h}| \frac{1 - 2^{-\mu} - 2.u}{1 - 2^{-\mu}}$$

So we have for $|\delta|$:

$$|\delta| \le |a_h + a_\ell + b_h + b_\ell| \cdot \frac{1 - 2^{-\mu}}{1 - 2^{-\mu} - 2 \cdot u} \cdot 3 \cdot u^2 \cdot \frac{1 - 2^{-\mu - 1}}{1 - 2^{-\mu}}$$
$$|\delta| \le |a_h + a_\ell + b_h + b_\ell| \cdot 3 \cdot u^2 \cdot \frac{1 - 2^{-\mu - 1}}{1 - 2^{-\mu} - 2 \cdot u}$$

So:

$$|\epsilon| \le 3.u^2 \cdot \frac{1 - 2^{-\mu - 1}}{1 - 2^{-\mu} - 2.u}$$

As $\mu \geq 1$, we want the upper bound of $\frac{1-2^{-\mu-1}}{1-2^{-\mu}-2.u} \Rightarrow$

$$\frac{3/4}{1/2 - 2.u} \le 2$$

So:

$$|\epsilon| \le 2.3.u^2 = 6.u^2$$

3.2 Multiplication Operators

3.2.1 Mul112

See algorithm 4

Lemma 4 (Mul112) Let a and b floating point numbers, s and t result of Mul112(a,b) for the 4 modes of rounding, considering that there is no **overflow** so:

- (1) $r_1 + r_2$ is exactly equal to a.b.
- $(2) |r_2| \leq u.|r_1|$

PROOF (1) According to [19], FMA(a, b, -c) is exactly equal to a + b - c for the 4 modes of rounding.

Based on Algorithm Mul112, we have :

 $r_2 = FMA(a, b, -r_1)$ that supposed $\Rightarrow r_2 = a * b - r_1$ So

$$r_1 + r_2 = a * b - r_1 + r_1$$

 $r_1 + r_2 = a * b$

(2) $r1 = \circ(a*b)$ according to the collary $6 \Rightarrow r1 = a*b + \epsilon$ with $|\epsilon| \le u.|r_1|$ We know that r2 is exactly equal to $a*b - r_1$. $\Rightarrow r_2 = a*b - r_1 \Rightarrow$

$$r_2 = a * b - (a * b + \epsilon)$$
$$r_2 = -\epsilon$$

That $|\epsilon| \leq u.|r_1| \Rightarrow |r_2| \leq u.|r_1|$.

3.2.2 Mul122

See algorithm 5

Lemma 5 (Mul122) Let a a **Double** number and (b_h, b_ℓ) a **Double-Double** number, r_1 and r_2 result of $Mul122(a, b_h, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so $|r_2| \le u.|r_1|$.

PROOF We know $|b_{\ell}| \leq u.|b_h|$

In the algorithme Mul122, we see $r_1, r_2 = Add112(t_1, t_6)$, bases on the lemma 1, we have $|r_2| \leq u.|r_1|$.

Theorem 5 (Relative error algorithm Mul122) Let a a double number and (b_h, b_ℓ) a double-double number are the arguments of the function Mul122:

So:

$$r_1 + r_2 = (a.(b_h + b_\ell)).(1 + \epsilon)$$

with $|\epsilon| \leq 3.u^2$.

 $3.u^2=2^{-104.41}$ for RNDN and $3.u^2=2^{-102.41}$ for the other rounding modes.

PROOF According to the calculation technique of the theorem 4.7 proof's([12]). We have from Algorithm 5:

$$r_1, r_2 = Add112(t_1, t_4)$$
 according to $Add112 \Rightarrow r_1 + r_2 = t_1 + t_4$
 $t_4 = \circ(t_2 + t_3) \Rightarrow t_4 = (t_2 + t_3)(1 + \epsilon_1)$ with $|\epsilon_1| \leq u$
 $t_3 = \circ(a.b_\ell) \Rightarrow t_3 = (a.b_\ell)(1 + \epsilon_2)$ with $|\epsilon_2| \leq u$
 $t_1, t_2 = Mul112(a, b_h)$ according to $Mul112 \Rightarrow |t_2| \leq u.|t_1|$ and $t_1 + t_2 = a \times b_h$

So we have:

$$t_4 = (t_2 + a.b_\ell.(1 + \epsilon_2)).(1 + \epsilon_1)$$

$$t_4 = t_2 + a.b_{\ell} + a.b_{\ell}.\epsilon_2 + t_2.\epsilon_1 + a.b_{\ell}.\epsilon_1 + a.b_{\ell}.\epsilon_2.\epsilon_1$$

as
$$r_1 + r_2 = t_1 + t_4 \Rightarrow$$

$$r_1 + r_2 = t_1 + (t_2 + a.b_{\ell} + a.b_{\ell}.\epsilon_2 + t_2.\epsilon_1 + a.b_{\ell}.\epsilon_1 + a.b_{\ell}.\epsilon_2.\epsilon_1)$$

but $t_1 + t_2 = a.b_h \Rightarrow$

$$r_1 + r_2 = a.b_h + a.b_\ell + a.b_\ell.\epsilon_2 + t_2.\epsilon_1 + a.b_\ell.\epsilon_1 + a.b_\ell.\epsilon_2.\epsilon_1$$

$$r_1 + r_2 = a.b_h + a.b_\ell + \delta$$

with

$$\delta = t_2 \cdot \epsilon_1 + a \cdot b_\ell \cdot \epsilon_2 + a \cdot b_\ell \cdot \epsilon_1 + a \cdot b_\ell \cdot \epsilon_2 \cdot \epsilon_1$$

So we have:

$$|\delta| \le |t_2.\epsilon_1| + |a.b_l.(\epsilon_1 + \epsilon_2 + \epsilon_2.\epsilon_1)|$$

as we know that $|\epsilon_1| \leq u$ and $|\epsilon_2| \leq u \Rightarrow$

$$|\delta| \le |t_2| \cdot u + |a \cdot b_\ell| \cdot (u^2 + 2 \cdot u)$$

According to the conditions of the algorithm:

$$|b_{\ell}| \leq u.|b_h| \Rightarrow |a \times b_{\ell}| \leq u.|a.b_h|$$
. So:

$$|\delta| \le |t_2| \cdot u + u \cdot |a \cdot b_h| \cdot (u^2 + 2 \cdot u)$$

$$|\delta| \le |t_2| \cdot u + |a \cdot b_h| \cdot (u^3 + 2 \cdot u^2)$$

As $a.b_h = t_1 + t_2$ and that $|t_2| \le u.|t_1|$, we have $|a.b_h| \ge |(1 + \frac{1}{u}).t_2|$ but $1 + \frac{1}{u} > 0 \Rightarrow$

$$(1 + \frac{1}{u}).|t_1| \le |a.b_h|$$

 $\frac{1+u}{u}.|t_2| \le |a.b_h|$
 $|t_2| \le \frac{u}{1+u}.|a.b_h|$

So:

$$|\delta| \le \frac{u^2}{1+u}.|a.b_h| + |a.b_h|.(u^3 + 2.u^2)$$

After calculate, \Rightarrow

$$|\delta| \le \frac{u^4 + 3.u^3 + 3.u^2}{1 + u}.|a.b_h|$$

Now, we search a lower bound for $|a_h.(b_h + b_\ell)|$ in fuction of $|a_h.b_h|$. For to calculate this lower bound, we need to search the upper bound of $|a.b_\ell|$. So we have:

$$|a.b_{\ell}| \leq u.|a.b_{h}|$$

So we have:

$$|a.(b_h + b_\ell)| \ge (1 - u).|a_h.b_h|$$

$$\frac{1}{1 - u}.|a.(b_h + b_\ell)| \ge |a_h.b_h|$$

For the calculation of relative error, we have:

$$|\delta| \le \frac{u^4 + 3.u^3 + 3.u^2}{1 - u^2}.|a.b_h|$$

As
$$\frac{u^4+3.u^3+3.u^2}{1-u^2} \le 3.u^2 \Rightarrow |\delta| \le 3.u^2.|a.(b_h+b_\ell)|.$$

For $r_1+r_2=(a.(b_h+b_\ell)).(1+\epsilon)$ We have that $|\epsilon| \le 3.u^2$

3.2.3 Mul222

See algorithm 6

Lemma 6 (Mul222) Let (a_h, a_ℓ) and (b_h, b_ℓ) are **Double-Double** numbers, r_1 and r_2 result of $Mul222(a_h, a_\ell, b_h, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so $|r_1| \le u.|r_2|$.

PROOF We know $|b_{\ell}| \leq u.|b_h|$ In the algorithme Mul222, we see $r_1, r_2 = Add112(t_1, t_6)$, based on the lemma 1, we have $|r_2| \leq u.|r_1|$.

Theorem 6 (Relative error algorithm Mul222) Let (a_h, a_ℓ) and (b_h, b_ℓ) are double-double numbers and are the arguments of the function Mul222: so:

$$r_1 + r_2 = ((a_h + a_\ell) \times (b_h + b_\ell))(1 + \epsilon)$$

with $|\epsilon| \leq 7.u^2$.

 $7.u^2=2^{-103.19}$ for RNDN and $7.u^2=2^{-101.19}$ for the other rounding modes.

PROOF According to the calculation technique of the theorem 4.7 proof's([12]). We have from Algorithm 7:

$$r_1, r_2 = Add112(t_1, t_6)$$
 according to $Add112 \Rightarrow r_1 + r_2 = t_1 + t_6$
 $t_6 = \circ(t_2 + t_5) \Rightarrow t_6 = (t_2 + t_5)(1 + \epsilon_1)$ with $|\epsilon_1| \leq u$
 $t_5 = \circ(t_3 + t_4) \Rightarrow t_5 = (t_3 + t_4)(1 + \epsilon_2)$ with $|\epsilon_2| \leq u$
 $t_4 = \circ(b_h \times a_\ell) \Rightarrow t_4 = (b_h.a_\ell).(1 + \epsilon_3)$ with $|\epsilon_3| \leq u$
 $t_3 = \circ(a_h \times b_\ell) \Rightarrow t_3 = (a_h.b_l).(1 + \epsilon_4)$ with $|\epsilon_4| \leq u$
 $t_1, t_2 = Mul112(a_h, b_h)$ according to $Mul112 \Rightarrow |t_2| \leq u.|t_1|$ and $t_1 + t_2 = a_h.b_h$

So we have:

$$t_6 = (t_2 + ((a_h.b_\ell).(1 + \epsilon_4) + (b_h.a_\ell).(1 + \epsilon_3)).(1 + \epsilon_2)).(1 + \epsilon_1)$$

As $r_1 + r_2 = t_1 + t_6 \Rightarrow$

$$r_1 + r_2 = t_1 + (t_2 + ((a_h.b_\ell).(1 + \epsilon_4) + (b_h.a_\ell).(1 + \epsilon_3)).(1 + \epsilon_2)).(1 + \epsilon_1)$$
$$r_1 + r_2 = t_1 + t_2 + a_h.b_\ell + b_h.a_\ell + \delta$$

with

$$\delta = t_2.\epsilon_1 + a_h.b_l.(\epsilon_4 + \epsilon_2 + \epsilon_4.\epsilon_2 + \epsilon_1 + \epsilon_4.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_4.\epsilon_2.\epsilon_1) + a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)$$

So we have:

$$|\delta| \leq |t_2.\epsilon_1| + |a_h.b_l.(\epsilon_4 + \epsilon_2 + \epsilon_4.\epsilon_2 + \epsilon_1 + \epsilon_4.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_4.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_3.\epsilon_1)| + |a_l.b_h.(\epsilon_3 + \epsilon_3.\epsilon_1)$$

$$|\delta| \leq |t_2.\epsilon_1| + |a_h.b_l|. |\epsilon_4 + \epsilon_2 + \epsilon_4.\epsilon_2 + \epsilon_1 + \epsilon_4.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_4.\epsilon_2.\epsilon_1| + |a_l.b_h|. |\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_3.\epsilon_2.$$

as we know that $|\epsilon_i| \leq u$ with $1 \leq i \leq 4 \Rightarrow$

$$|\delta| \le |t_2.\epsilon_1| + (|a_h.b_l| + |a_l.b_h|).|\epsilon_4 + \epsilon_2 + \epsilon_4.\epsilon_2 + \epsilon_1 + \epsilon_4.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_4.\epsilon_2.\epsilon_1|$$

According to the conditions of the algorithm:

 $|a_\ell| \le u.|a_h|$ and $|b_\ell| \le u.|b_h| \Rightarrow |a_\ell \times b_h| \le u.|a_h \times b_h|$ and $|a_h \times b_\ell| \le u.|a_h \times b_h|$ $u.|a_h \times b_h|$. So:

$$|\delta| \le |t_2.\epsilon_1| + (u.|a_h.b_h| + u|a_h.b_h|).|\epsilon_4 + \epsilon_2 + \epsilon_4.\epsilon_2 + \epsilon_1 + \epsilon_4.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_4.\epsilon_2.\epsilon_1|$$

We search the upper bound of $|t_2|$ in function of $|a_h.b_h|$, we begin with the upper bound of $|t_1|$, as $t_1 + t_2 = a_h b_h$ and $|t_2| \le u |t_1|$, we have that:

$$|t_1 + t_2| \ge (1 - u).|t_1|$$

 $|t_1| \le \frac{1}{1 - u}.|a_h.b_h|$

and so:

$$|t_2| \le \frac{u}{1-u}.|a_h.b_h|$$

$$|\delta| \le \frac{u}{1-u}.|a_h.b_h|.|\epsilon_1| + 2.u.|a_h.b_h|.|\epsilon_4 + \epsilon_2 + \epsilon_4.\epsilon_2 + \epsilon_1 + \epsilon_4.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_4.\epsilon_2.\epsilon_1|$$

$$|\delta| \le |a_h.b_h|.(\frac{u}{1-u}.|\epsilon_1| + 2.u.|\epsilon_4 + \epsilon_2 + \epsilon_4.\epsilon_2 + \epsilon_1 + \epsilon_4.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_4.\epsilon_2.\epsilon_1|)$$

After calculation, \Rightarrow

$$|\delta| \le |a_h.b_h|.\frac{-2.u^5 - 4.u^4 + 7.u^2}{1 - u}$$

Now, we search a lower bound for $|(a_h + a_\ell).(b_h + b_\ell)|$ in function of $|a_h.b_h|$. To calculate this lower bound, we need to search the upper bound of $|a_h.b_\ell+$ $a_{\ell}.b_h + a_{\ell}.b_{\ell}$. So we have:

$$|a_h.b_{\ell} + a_{\ell}.b_h + a_{\ell}.b_{\ell}| \le |a_h.b_{\ell}| + |a_{\ell}.b_h| + |a_{\ell}.b_{\ell}|$$

$$|a_h.b_{\ell} + a_{\ell}.b_h + a_{\ell}.b_{\ell}| \le u.|a_h.b_h| + u.|a_h.b_h| + u^2.|a_h.b_h|$$

$$|a_h.b_{\ell} + a_{\ell}.b_h + a_{\ell}.b_{\ell}| \le (u^2 + 2.u).|a_h.b_h|$$

So we have:

$$|(a_h + a_\ell).(b_h + b_\ell)| \ge (1 - (u^2 + 2.u)).|a_h.b_h|$$

For the calculation of relative error, we have:

$$|\delta| \le \frac{-2 \cdot u^5 - 4 \cdot u^4 + 7 \cdot u^2}{(1 - u) \cdot (1 - (u^2 + 2 \cdot u))} \cdot |(a_h + a_\ell) \cdot (b_h + b_\ell)|$$

but
$$\frac{-2.u^5 - 4.u^4 + 7.u^2}{(1-u).(1-(u^2+2.u))} \le 7.u^2 \Rightarrow |\delta| \le 7.u^2.|(a_h + a_\ell).(b_h + b_\ell)|.$$

For $r_1 + r_2 = ((a_h + a_\ell).(b_h + b_\ell)).(1 + \epsilon)$ so $|\epsilon| \le 7.u^2$

Chapter 4

Operators on Triple-Double numbers

4.1 Addition Operators

4.1.1 Add133

See algorithm 7

Lemma 7 (Add133) Let a a **double** number and (b_h, b_m, b_ℓ) a **triple-Double** number, r_h , r_m and r_ℓ result of $Add133(a, b_h, b_m, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so $|r_\ell| \le u.|r_m|$, $|r_m| \le u.|r_h|$ and $|r_\ell| \le u^2.|r_h|$

PROOF We suppose that $|b_m| \le u.|b_h|$, $|b_\ell| \le u.|b_m|$, according to Add133, so we have:

- $r_m + r_\ell = Add112(t_2, t_4)$ thanks to $Add112 \Rightarrow r_m + r_\ell = t_2 + t_4$ and $|r_\ell| \leq u.|r_m|$.
- $t_4 = \circ(t_3 + b_\ell) \Rightarrow t_4 = (t_3 + b_\ell).(1 + \epsilon)$ with $|\epsilon| \le u$.
- $t_2 + t_3 = Add112(t_1, b_m)$ thanks to $Add112 \Rightarrow t_2 + t_3 = t_1 + b_m$ and $|t_3| \leq u.|t_2|$.
- $r_h + t_1 = Add112(a, b_h)$ thanks to $Add112 \Rightarrow r_h + t_1 = a + b_h$ and $|t_1| \leq u.|r_h| \Rightarrow$

$$r_m + r_\ell = t_2 + t_4$$

$$r_m + r_\ell = t_2 + (t_3 + b_\ell).(1 + \epsilon)$$

$$r_m + r_\ell = t_2 + t_3 \cdot (1 + \epsilon) + b_\ell \cdot (1 + \epsilon)$$

 $r_m + r_\ell = t_2 + t_3 + t_3 \cdot \epsilon + b_\ell \cdot (1 + \epsilon)$

As $t_2 + t_3 = t_1 + b_m \Rightarrow$

$$r_m + r_{\ell} = t_1 + b_m + t_3 \cdot \epsilon + b_{\ell} \cdot (1 + \epsilon)$$
$$|r_m + r_{\ell}| = |t_1 + b_m + t_3 \cdot \epsilon + b_{\ell} \cdot (1 + \epsilon)|$$
$$|r_m + r_{\ell}| \le |t_1 + b_m| + |t_3 \cdot \epsilon| + |b_{\ell} \cdot (1 + \epsilon)|$$

We search the lower bound of $|t_3|$ in function of $|t_1+b_m|$: $|t_1+b_m|=|t_2+t_3|\Rightarrow |t_1+b_m|\geq (1+\frac{1}{u}).|t_3|\Rightarrow |t_3|\leq \frac{u}{u+1}.|t_1+b_m|$

$$|r_m + r_\ell| \le |t_1 + b_m| + \frac{u}{u+1} \cdot |(t_1 + b_m) \cdot \epsilon| + |b_\ell \cdot (1+\epsilon)|$$

$$|r_m + r_\ell| \le |t_1 + b_m| + \frac{u^2}{u+1} \cdot |(t_1 + b_m)| + |b_\ell \cdot (1+u)|$$

$$|r_m + r_\ell| \le +\frac{u^2 + u + 1}{u + 1}.|(t_1 + b_m)| + |b_\ell.(1 + u)|$$

According to the conditions of Add112, we have $|t_1| \ge |b_m|$ and also that $|b_\ell| \le u.|b_m| \Rightarrow$

$$|r_m + r_\ell| \le \frac{2.(u^2 + u + 1) + u(u + 1)^2}{u + 1}.|t_1|$$

 $|r_m + r_\ell| \le \frac{u^3 + 4.u^2 + 3.u + 1}{u + 1}.|t_1|$

As $|t_1| \leq u.|r_h| \Rightarrow$

$$|r_m + r_\ell| \le u \cdot \frac{u^3 + 4 \cdot u^2 + 3 \cdot u + 1}{u + 1} \cdot |r_h|$$

We have: $\frac{u^4+4.u^3+3.u^2+u}{u+1} \le u \Rightarrow$

$$|r_m + r_\ell| \le u.|r_h|$$

We know that $|r_{\ell}| \leq u.|r_m| \Rightarrow (1-u).|r_m| \leq |r_m + r_{\ell}|$.

We have:

$$(1-u).|r_m| \le u.|r_h|$$
$$|r_m| \le \frac{u}{1-u}.|r_h|$$

But $\frac{u}{1-u} \le u \Rightarrow$

$$|r_m| \le u.|r_h|$$

and

$$|r_{\ell}| \le u^2.|r_h|$$

47

Theorem 7 (Relative error algorithm Add133 without occurring of cancellation) Let a a double number and (b_h, b_m, b_ℓ) a triple-double number are the arguments of the function Add133. So:

$$r_h + r_m + r_\ell = (a + (b_h + b_m + b_\ell)).(1 + \epsilon)$$

with $|\epsilon| \leq 3.u^3$.

 $3.u^3=2^{-157.41}$ for RNDN and $3.u^3=2^{-154.41}$ for the other rounding modes.

PROOF According to the calculation technique of the theorem 4.2 proof's([12]). We have $|a| \ge |b_h|$, either a > 0 and $b_h > 0$ or a < 0 and $b_h < 0$. As they're symmetric, we only use a > 0 and $b_h > 0$.

We have from Algorithme 7:

 $r_m + r_\ell = Add112(t_2, t_4)$ based on $Add112 \Rightarrow r_m + r_\ell = t_2 + t_4$ $t_4 = \circ(t_3 + b_\ell) \Rightarrow t_4 = (t_3 + b_\ell)(1 + \epsilon_1)$ with $|\epsilon_1| \leq u$ $t_2 + t_3 = Add112(t_1, b_m)$ based on $Add112 \Rightarrow t_2 + t_3 = t_1 + b_m$ $r_h + t_1 = Add112(a, b_h)$ based on $Add112 \Rightarrow r_h + t_1 = a + b_h$ So we have:

$$r_m + r_{\ell} = t_2 + (t_3 + b_{\ell})(1 + \epsilon_1)$$
$$r_m + r_{\ell} = t_2 + t_3 + t_3 \cdot \epsilon + b_{\ell} \cdot (1 + \epsilon_1)$$

As $t_2 + t_3 = t_1 + b_m \Rightarrow$

$$r_m + r_\ell = t_1 + b_m + t_3 \cdot \epsilon_1 + b_\ell \cdot (1 + \epsilon_1)$$

Then we have:

$$r_h + r_m + r_\ell = r_h + t_1 + b_m + t_3 \cdot \epsilon_1 + b_\ell \cdot (1 + \epsilon_1)$$

As $r_h + t_1 = a + b_h$

$$r_h + r_m + r_\ell = a + b_h + b_m + t_3 \cdot \epsilon_1 + b_\ell \cdot (1 + \epsilon_1)$$

$$r_h + r_m + r_\ell = a + (b_h + b_m + b_\ell + t_3 \cdot \epsilon_1 + b_\ell \cdot \epsilon_1$$

$$r_h + r_m + r_\ell = a + (b_h + b_m + b_\ell) + \delta$$

with $\delta = t_3 \cdot \epsilon_1 + b_\ell \cdot \epsilon_1$

So:

$$|\delta| \le |t_3.\epsilon_1| + |b_\ell.\epsilon_1|$$

We know that $t_1 + b_m = t_2 + t_3$ and that $|t_3| \le u.|t_2| \Rightarrow$

$$|t_1 + b_m| \ge \left(\frac{1}{u} + 1\right).|t_3|$$

$$|t_3| \le \frac{u}{u+1}.|t_1 + b_m|$$

$$|t_3| \le \frac{u}{u+1}.(|t_1| + |b_m|)$$

As $a + b_h = r_h + t_1$ and that $|t_1| \le u.|r_h| \Rightarrow$

$$|a + b_h| \ge (\frac{1}{u} + 1).|t_1|$$

 $|t_1| \le \frac{u}{u+1}.|a + b_h|$

 \Rightarrow

$$|\delta| \le \frac{u}{u+1} \cdot (\frac{u}{u+1} \cdot |a+b_h| + |b_m|) \cdot \epsilon_1 | + |b_\ell \cdot \epsilon_1|$$

As $|\epsilon_1| \le u$, $|b_m| \le u |b_h| \le u |a+b_h|$ and $|b_\ell| \le u |b_m| \Rightarrow$

$$|\delta| \le \frac{u}{u+1} \cdot (\frac{u}{u+1} \cdot |a+b_h| + u \cdot |a+b_h|) \cdot u + u \cdot u^2 \cdot |a+b_h|$$

After calculating, we find:

$$|\delta| \le \frac{u^5 + 3.u^4 + 3.u^3}{(u+1)^2}.|a+b_h|$$

We search the upper bound of $|a + b_h|$ in function of $|a + b_h| + b_m + b_\ell|$: first, we calculate the upper bound of $|b_m + b_\ell|$ in function of $|a + b_h|$:

$$|b_m + b_\ell| \le |b_m| + |b_\ell|$$

 $|b_m + b_\ell| \le u.|a + b_h| + u^2.|a + b_h|$

We deduce:

$$|a + b_h + b_m + b_\ell| \ge (1 - (u + u^2)).|a + b_h|$$

 $|a + b_h + b_m + b_\ell| \ge |(1 - u - u^2)|.|a + b_h|$

$$\frac{1}{(1-u-u^2)} |a+b_h+b_m+b_\ell| \ge |a|$$

 \Rightarrow

$$|\delta| \le \frac{u^5 + 3 \cdot u^4 + 3 \cdot u^3}{(u+1)^2 \cdot (1-u-u^2)} \cdot |a+b_h+b_m+b_\ell|$$
$$|\delta| \le |\epsilon| \cdot |a+b_h+b_m+b_\ell|$$

We have $\frac{u^5+3.u^4+3.u^3}{(u+1)^2.(1-u-u^2)} \le 3.u^3 \Rightarrow |\epsilon| \le 3.u^3$

49

Theorem 8 (Relative error algorithm Add133 with a bounded cancellation)

Let a a **double** number and (b_h, b_m, b_ℓ) a **triple-double** number. We have for the algorithm Add133 with a and (b_h, b_m, b_ℓ) it's arguments. If a and b_h are different sign and we suppose $|b_h| \leq 2^{-\mu} |a|$ with $\mu \geq 1$. So:

$$r_h + r_m + r_\ell = (a + (b_h + b_m + b_\ell)).(1 + \epsilon)$$

with

$$|\epsilon| \le \frac{u^5 + 3.u^4 + 3.u^3}{(u+1)^2} \cdot \frac{1 - 2^{-\mu}}{1 - (1 + u + u^2) \cdot 2^{-\mu}} \le 6.u^3$$

PROOF According to the calculation technique of the theorem 4.3 proof's([12]). We suppose $|b_h| \leq 2^{-\mu}|a|$ with $\mu \geq 1$. First we look for the upper bound, using the results of the proof of Theorem 7. We have:

$$|b_{\ell}| \le u.|b_m| \le u^2.|b_h| \le u^2.2^{-\mu}|a|$$

then we search the lower bound of $|a + b_h|$ in function of |a|:

$$|a+b_h| \ge (1-2^{-\mu}).|a|$$

Now, we search the upper bound of $|b_m|$ and $|b_h|$ in function of $|a + b_h|$.

$$|b_m| \le \frac{u \cdot 2^{-\mu}}{1 - 2^{-\mu}} |a + b_h|$$

and

$$|b_{\ell}| \le \frac{u^2 \cdot 2^{-\mu}}{1 - 2^{-\mu}} |a + b_h|$$

As:

$$|\delta| \le \frac{u^5 + 3.u^4 + 3.u^3}{(u+1)^2}.|a+b_h|$$

Now, we search the lower bound for $|a + b_h + b_m + b_\ell|$ depending on $|a + b_h|$. We have:

$$|b_m + b_{\ell}| \le |b_m| + |b_{\ell}|$$

$$|b_m + b_{\ell}| \le (u + u^2) \cdot 2^{-\mu} \cdot |a|$$

$$|b_m + b_{\ell}| \le \frac{(u + u^2) \cdot 2^{-\mu}}{1 - 2^{-\mu}} |a + b_h|$$

$$|a + b_h + b_m + b_{\ell}| \ge |a_h + b_h| \frac{1 - 2^{-\mu} - (u + u^2) \cdot 2^{-\mu}}{1 - 2^{-\mu}}$$

So we have for $|\delta|$:

$$|\delta| \le |a+b_h+b_m+b_\ell| \cdot \frac{u^5 + 3 \cdot u^4 + 3 \cdot u^3}{(u+1)^2} \cdot \frac{1 - 2^{-\mu}}{1 - (1 + u + u^2) \cdot 2^{-\mu}}$$

So:

$$|\epsilon| \le \frac{u^5 + 3.u^4 + 3.u^3}{(u+1)^2} \cdot \frac{1 - 2^{-\mu}}{1 - (1 + u + u^2) \cdot 2^{-\mu}}$$

As $\mu \geq 1$, we want the upper bound of $\frac{1-2^{-\mu}}{1-(1+u+u^2)\cdot 2^{-\mu}} \Rightarrow$

$$\frac{1}{1/2} \le 2$$

So:

$$|\epsilon| \le 6.u^3$$

4.1.2 Add333

See algorithm 8

Lemma 8 (Add333) Let (a_h, a_m, a_ℓ) and (b_h, b_m, b_ℓ) are **triple-Double** number, r_h , r_m and r_ℓ result of $Add333(a_h, a_m, a_\ell, b_h, b_m, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so: $|r_\ell| \le u.|r_m|, |r_m| \le 2^{-2}.|r_h|$ and $|r_\ell| \le 2^{-2}.u.|r_h|$

PROOF • $r_h + r_\ell = Add112(t_7, t_8)$ According to $Add112 \Rightarrow |r_\ell| \leq u.|r_h|$.

- $t_8 = \circ(t_5 + t_6)$ according to collary 5: $t_8 = (t_5 + t_6).(1 + \epsilon_1)$ with $|\epsilon_1| \leq u$
- $t_5 = \circ(t_3 + t_4)$ according to collary 5: $t_5 = (t_3 + t_4).(1 + \epsilon_2)$ with $|\epsilon_2| \le u$
- $t_6 = \circ(a_\ell + b_\ell)$ according to collary 5: $t_6 = (a_\ell + b_\ell).(1 + \epsilon_3)$ with $|\epsilon_3| \leq u$.
- $t_7 + t_4 = Add112(t_1, t_2)$ according to Add112: $|t_4| \le u.|t_7| \Rightarrow |t_7 + t_4| \ge (1-u).|t_7| \Rightarrow (1-u).|t_7| \le |t_1 + t_2| \Rightarrow |t_7| \le \frac{1}{1-u}.|t_1 + t_2| \Rightarrow |t_4| \le \frac{u}{1-u}.|t_1 + t_2|$
- $t_2+t_3=Add112Cond(a_m,b_m)$ according to Add112: $|t_3|\leq u.|t_2|\Rightarrow |t_3|\leq \frac{u}{1-u}.|a_m+b_m|$

• $|t_5| \le (|t_3| + |t_4|)(1+u) \Rightarrow |t_5| \le (\frac{u}{1-u}.|a_m + b_m| + \frac{u}{1-u}.|t_1 + t_2|)(1+u)$ thanks to the condition $Add112 \Rightarrow |t_1| \ge |t_2| \Rightarrow$

$$|t_5| \le \left(\frac{u}{1-u}.|a_m + b_m| + \frac{2.u}{1-u}.|t_1|\right)(1+u)$$

$$|t_5| \le \left(u.|a_m + b_m| + 2.u.|t_1|\right)$$

As $|t_1| < u.|r_h| \Rightarrow$

$$|t_5| \le (u.|a_m + b_m| + 2.u^2.|r_h|)$$

As $|a_m| \le 2^{-\alpha_o} . |a_h|$ and that $|b_m| \le 2^{-\beta_o} . |b_h| \le 2^{-\beta_o} . \frac{3}{4} . |a_h|$, so we have:

$$|a_m + b_m| \le (2^{-\alpha_o} + 2^{-\beta_o} \cdot \frac{3}{4}).|a_h|$$

We search the lower bound of $|a_h|$ in function of $|a_h+b_h|$ and as $b_h \leq \frac{3}{4}$. $|a_h|$

$$|a_h + b_h| \ge (1 - \frac{3}{4}).|a_h|$$

 $|a_h| \le 4.|a_h + b_h|$

 \Rightarrow

$$|a_m + b_m| \le (2^{-\alpha_o} + 2^{-\beta_o} \cdot \frac{3}{4}) \cdot 4 \cdot |a_h + b_h|$$

 $|a_m + b_m| \le (2^{-\alpha_o + 2} + 2^{-\beta_o} \cdot 3) \cdot |a_h + b_h|$

But we know that $|a_h + b_h| \le (1 + u).|r_h|$

$$|a_m + b_m| \le (2^{-\alpha_o + 2} + 2^{-\beta_o} \cdot 3) \cdot (1 + u) \cdot |r_h|$$

but $\alpha_o \ge 4$ and $\beta_o \ge 4$

$$|a_m + b_m| \le 7.2^{-4} \cdot (1+u) \cdot |r_h|$$

$$|t_5| \le (u.7.2^{-4}.(1+u).|r_h| + 2.u^2.|r_h|)$$

After calculating, we have

$$|t_5| \le (2 + 7.2^{-4}).u^2 + 7.2^{-4}.u).|r_h|$$

We know that $t_6 = (a_{\ell} + b_{\ell})(1 + \epsilon_3) \Rightarrow$

$$|t_6| \le |a_\ell + b_\ell| \cdot (1+u)$$

As $|a_{\ell}| \le 2^{-\alpha_u} . |a_m| \le 2^{-\alpha_u - \alpha_o} . |a_h|$ and that $|b_{\ell}| \le 2^{-\beta_u} . |b_m| \le 2^{-\beta_u - \beta_o} . |b_h| \le 2^{-\beta_u - \beta_o} . \frac{3}{4} . |a_h|$, so we have:

$$|a_{\ell} + b_{\ell}| \le (2^{-\alpha_u - \alpha_o} + 2^{-\beta_u - \beta_o} \cdot \frac{3}{4}).|a_h|$$

$$|a_{\ell} + b_{\ell}| < (2^{-\alpha_u - \alpha_o + 2} + 2^{-\beta_u - \beta_o}.3).|a_h + b_h|$$

As $\alpha_o \geq 4$, $\alpha_u \geq 1$, $\beta_o \geq 4$ and $\beta_u \geq 1$, \Rightarrow

$$|a_{\ell} + b_{\ell}| < (2^{-3} + 2^{-5}.3).|a_h + b_h|$$

$$|a_{\ell} + b_{\ell}| \le 7.2^{-5}.|a_h + b_h|$$

 \Rightarrow

$$|t_6| \le 7.2^{-5} \cdot (1+u) \cdot |r_h| \cdot (1+u)$$

$$|t_6| \le (7.2^{-5} + 7.2^{-4}.u + 7.2^{-5}.u^2).|r_h|$$

 \Rightarrow

$$|t_8| \le (|t_5| + |t_6|).(1+u)$$

$$|t_8| \le ((2+7.2^{-4}).u^2+7.2^{-4}.u).|r_h| + (7.2^{-5}+7.2^{-4}.u+7.2^{-5}.u^2).|r_h|).(1+u)$$

After the calculations, we have :

$$|t_8| \leq ((2+7.2^{-4}+7.2^{-5}).u^3 + (2+7.2^{-4}+7.2^{-5}+7.2^{-3}).u^2 + (7.2^{-5}+7.2^{-3}).u + 7.2^{-5}).|r_h|$$

Now we look at the calculation of upper bound of $|t_7|$ in function of $|r_h|$. we have $|t_7| \leq \frac{1}{1-u}.|t_1+t_2|$ thanks to the condition of Add112, we have:

$$|t_7| \leq frac21 - u.|t_1|$$

As $|t_1| \leq u.|r_h| \Rightarrow$

$$|t_7| \le \frac{2.u}{1-u}.|r_h|$$

We have:

$$|r_m + r_\ell| \le |t_7| + |t_8|$$

$$|r_m + r_\ell| \le \frac{2.u}{1 - u}.|r_h| + ((2 + 7.2^{-4} + 7.2^{-5}).u^3 + (2 + 7.2^{-4} + 7.2^{-5} + 7.2^{-3}).u^2 + (7.2^{-5} + 7.2^{-3}).u + 7.2^{-5} + 7.2^{-5}).u^3 + (2 + 7.2^{-4} + 7.2^{-5} + 7.2^{-3}).u^2 + (7.2^{-5} + 7.2^{-3}).u + 7.2^{-5} + 7.2^{-5} + 7.2^{-5}).u^3 + (2 + 7.2^{-4} + 7.2^{-5} + 7.2^{-3}).u^3 + (2 + 7.2^{-5} + 7.2^{-5}).u^3 + (2 + 7.2^$$

We have
$$\frac{2.u}{1-u} + ((2+7.2^{-4}+7.2^{-5}).u^3 + (2+7.2^{-4}+7.2^{-5}+7.2^{-3}).u^2 + (7.2^{-5}+7.2^{-3}).u + 7.2^{-5}) \le 2^{-2} \Rightarrow$$

$$|r_m + r_\ell| \le 2^{-2}.|r_h|$$

4.1. ADDITION OPERATORS

53

As
$$|r_h| \le u.|r_\ell| \Rightarrow (1-u).|r_m| \le |r_m+r_\ell| \Rightarrow$$

$$|r_m| \le \frac{1}{1-u}.2^{-2}.|r_h|$$

$$|r_m| \le 2^{-2}.|r_h|$$
 and
$$|r_\ell| \le 2^{-2}.u.|r_h|$$

Theorem 9 (Relative error algorithm Add333) Let (a_h, a_m, a_ℓ) and (b_h, b_m, b_ℓ) are **triple-double** numbers and arguments of the function Add333. So:

$$r_h + r_m + r_\ell = ((a_h + a_m + a_\ell) + (b_h + b_m + b_\ell)).(1 + \epsilon)$$

$$with \ |\epsilon| \le (4.u^4 + 5.u^3 + u^2).2^{-min(\alpha_o,\beta_o)+5} + (u^2 + 2.u).2^{-min(\alpha_u + \alpha_o,\beta_u + \beta_o)+5}$$

PROOF According to the calculation technique of the theorem 5.2 prof's([12]). According to Add333 and thanks to Add112, we have these results:

$$|r_{\ell}| \le u.|r_{\ell}|$$

 $|t_4| \le u.|t_7|$
 $|t_3| \le u.|t_2|$
 $|t_1| \le u.|r_h|$

We search the upper bound of $|r_h|$ in function of $|a_h|$.

$$|t_1 + r_h| = |a_h + b_h|$$

As
$$|b_h| \le \frac{3}{4}.|a_h|$$

$$|t_1 + r_h| \le (1 + \frac{3}{4}).|a_h|$$

$$|t_1 + r_h| \le \frac{7}{4}.|a_h| \le 2.|a_h|$$

As
$$|t_1| \le u.|r_h|$$
 $|t_1 + r_h| \ge (1 - u).|r_h|$

So we have:

$$|r_h| \le \frac{2}{1-u}.|a_h|$$

As $|t_1| \leq u.|r_h|$

$$|t_1| \le \frac{2.u}{1-u}.|a_h|$$

Now, we calculate for all $|t_i|$ with $2 \le i \le 8$ their upper bounds in function of $|a_h|$. As calculated previously with r_h , we have:

$$|t_2| \le \frac{1}{1-u} . |a_m + b_m|$$

As $|a_m| \le 2^{-\alpha_o} . |a_h|$ and that $|b_m| \le 2^{-\beta_o} . |b_h| \le 2^{-\beta_o} . \frac{3}{4} . |a_h| \Rightarrow$

$$|t_2| \le \frac{1}{1-u} \cdot (2^{-\alpha_o} + 2^{-\beta_o} \cdot \frac{3}{4}) \cdot |a_h|$$

As $\frac{3}{4} \le 1 \Rightarrow$

$$|t_2| \le \frac{1}{1-u} \cdot (2^{-\alpha_o} + 2^{-\beta_o}) \cdot |a_h|$$

$$|t_2| \le \frac{1}{1-u} \cdot 2^{-\min(\alpha_o, \beta_o) + 1} \cdot |a_h|$$

As $|t_3| \le u.|t_2| \Rightarrow$

$$|t_3| \le \frac{u}{1-u} \cdot 2^{-\min(\alpha_o, \beta_o) + 1} \cdot |a_h|$$

Same for $|t_7|$, we have:

$$|t_7| \le \frac{1}{1-u}.|t_1 + t_2|$$

$$|t_7| \le \frac{1}{1-u}.(\frac{2.u}{1-u}.|a_h| + \frac{1}{1-u}.2^{-\min(\alpha_o,\beta_o)+1}.|a_h|)$$

$$|t_7| \le \frac{1}{(1-u)^2}.(2.u + 2^{-\min(\alpha_o,\beta_o)+1}).|a_h|$$

In case where $2.u = 2^{-min(\alpha_o,\beta_o)+1}$, we have:

$$|t_7| \le \frac{1}{(1-u)^2} \cdot 2^{-\min(\alpha_o, \beta_o) + 2} \cdot |a_h|$$

As $|t_4| \leq u.|t_7| \Rightarrow$

$$|t_4| \le \frac{u}{(1-u)^2} \cdot 2^{-\min(\alpha_o, \beta_o) + 2} \cdot |a_h|$$

For t_6 :

$$|t_6| \le (|a_\ell| + |b_\ell|).(1+u)$$

As
$$|a_{\ell}| \le 2^{-\alpha_u} . |a_m| \le 2^{-\alpha_u - \alpha_o} . |a_h|$$
 and $|b_{\ell}| \le 2^{-\beta_u} . |b_m| \le 2^{-\beta_u - \beta_o} . |b_h| \le 2^{-\beta_u - \beta_o} . \frac{3}{4} . |a_h| \Rightarrow$

$$|t_6| \le (2^{-\alpha_u - \alpha_o}.|a_h| + 2^{-\beta_u - \beta_o}.\frac{3}{4}.|a_h|).(1+u)$$

$$|t_6| \le (1+u) \cdot 2^{-\min(\alpha_u + \alpha_o, \beta_u + \beta_o) + 1} \cdot |a_h|$$

And t_5 , we have:

$$|t_{5}| \leq (|t_{3}| + |t_{4}|).(1+u)$$

$$|t_{5}| \leq (1+u).(\frac{u}{1-u}.2^{-\min(\alpha_{o},\beta_{o})+1}.|a_{h}| + \frac{u}{(1-u)^{2}}.2^{-\min(\alpha_{o},\beta_{o})+2}.|a_{h}|)$$

$$|t_{5}| \leq (1+u).\frac{u}{(1-u)^{2}}((1-u).2^{-\min(\alpha_{o},\beta_{o})+1}. + u.2^{-\min(\alpha_{o},\beta_{o})+2}).|a_{h}|$$

$$|t_{5}| \leq (1+u).\frac{u^{2}}{(1-u)^{2}}.(1+u).2^{-\min(\alpha_{o},\beta_{o})+1}.|a_{h}|$$

After the calculation:

$$|t_5| \le \frac{u^3 + 2.u^2 + u}{(1 - u)^2} \cdot 2^{-\min(\alpha_o, \beta_o) + 1} \cdot |a_h|$$

Simplifying, we have:
As
$$\frac{u^3+2.u^2+u}{(1-u)^2} \le u \Rightarrow$$

$$|t_5| \le u.2^{-\min(\alpha_o,\beta_o)+1}.|a_h|$$

and similarly for the others:

$$|r_h| \le 2.|a_h|$$

$$|t_1| \le 2.u.|a_h|$$

$$|t_2| \le 2^{-\min(\alpha_o,\beta_o)+1}.|a_h|$$

$$|t_3| \le u.2^{-\min(\alpha_o,\beta_o)+1}.|a_h|$$

$$|t_4| \le u^2.2^{-\min(\alpha_o,\beta_o)+2}.|a_h|$$

$$|t_6| \le 2^{-\min(\alpha_u+\alpha_o,\beta_u+\beta_o)+1}.|a_h|$$

$$|t_7| \le u.2^{-\min(\alpha_o,\beta_o)+2}.|a_h|$$

As $r_m + r_\ell = t_7 + t_8$, we calculated t_8 :

$$t_8 = (t_5 + t_6).(1 + \epsilon_1)$$

with $|\epsilon_1| \leq u$

$$t_8 = ((t_3 + t_4).(1 + \epsilon_2) + (a_\ell + b_\ell).(1 + \epsilon_3)).(1 + \epsilon_1)$$

with $|\epsilon_2| \leq u$ and $|\epsilon_3| \leq u$ After the calculations:

$$t_8 = t_3 + t_4 + a_\ell + b_\ell + \delta$$

with
$$\delta = t_3.(\epsilon_2 + \epsilon_1 + \epsilon_2.\epsilon_1) + t_4.(\epsilon_2 + \epsilon_1 + \epsilon_2.\epsilon_1) + a_\ell.(\epsilon_3 + \epsilon_1 + \epsilon_3.\epsilon_1) + b_\ell.(\epsilon_3 + \epsilon_1 + \epsilon_3.\epsilon_1)$$

We seek the upper bound of $|\delta|$ in function of $|a_h|$:

$$|\delta| \le |t_3 \cdot (\epsilon_2 + \epsilon_1 + \epsilon_2 \cdot \epsilon_1)| + |t_4 \cdot (\epsilon_2 + \epsilon_1 + \epsilon_2 \cdot \epsilon_1)| + |a_\ell \cdot (\epsilon_3 + \epsilon_1 + \epsilon_3 \cdot \epsilon_1)| + |b_\ell \cdot (\epsilon_3 + \epsilon_1 + \epsilon_3 \cdot \epsilon_1)|$$

As $|\epsilon_i| \le u$ with $1 \le i \le 3 \Rightarrow$

$$|\delta| < (2.u + u^2).(|t_3| + |t_4| + |a_{\ell}| + |b_{\ell}|)$$

$$|\delta| \le (2.u + u^2) \cdot (u \cdot 2^{-\min(\alpha_o, \beta_o) + 1} \cdot |a_h| + u^2 \cdot 2^{-\min(\alpha_o, \beta_o) + 2} \cdot |a_h| + 2^{-\alpha_u - \alpha_o} \cdot |a_h| + 2^{-\beta_u - \beta_o} \cdot \frac{3}{4} \cdot |a_h|)$$

$$|\delta| \leq (2.u + u^2).(u.2^{-min(\alpha_o,\beta_o)+1}.|a_h| + u^2.2^{-min(\alpha_o,\beta_o)+2}.|a_h| + 2^{-\alpha_u-\alpha_o}.|a_h| + 2^{-\beta_u-\beta_o}|a_h|)$$

$$|\delta| \le (2.u + u^2) \cdot (u \cdot 2^{-\min(\alpha_o, \beta_o) + 1} \cdot |a_h| + u^2 \cdot 2^{-\min(\alpha_o, \beta_o) + 2} \cdot |a_h| + 2^{-\min(\alpha_u + \alpha_o, \beta_u + \beta_o) + 1} \cdot |a_h|)$$

$$|\delta| \le ((4.u^4 + 5.u^3 + u^2).2^{-min(\alpha_o, \beta_o) + 1}. + (u^2 + 2.u).2^{-min(\alpha_u + \alpha_o, \beta_u + \beta_o) + 1}).|a_h|$$

We look the lower bound of $|a_h|$ in function of $|a_h + b_h + a_m + b_m + a_\ell + b_\ell|$:

$$|b_h + a_m + b_m + a_\ell + b_\ell| \le |b_h| + |a_m| + |b_m| + |a_\ell| + |b_\ell|$$

$$|b_h + a_m + b_m + a_\ell + b_\ell| \le \frac{3}{4}.|a_h| + 2^{-\alpha_o}.|a_h| + \frac{3}{4}.2^{-\beta_o}.|a_h| + 2^{-\alpha_o - \alpha_u}|a_h| + \frac{3}{4}.2^{-\beta_o - \beta_u}.|a_h|$$

with the condition of this algorithm, after calculate:

$$|b_h + a_m + b_m + a_\ell + b_\ell| \le \frac{117}{128} |a_h|$$

As $\frac{117}{128} < \frac{15}{16} \Rightarrow$

$$|b_h + a_m + b_m + a_\ell + b_\ell| \le \frac{15}{16}.|a_h|$$

We can say that:

$$|a_h + b_h + a_m + b_m + a_\ell + b_\ell| \ge \frac{1}{16}.|a_h|$$

We have:

$$|\delta| \leq 16.((4.u^4 + 5.u^3 + u^2).2^{-min(\alpha_o,\beta_o) + 1} + (u^2 + 2.u).2^{-min(\alpha_u + \alpha_o,\beta_u + \beta_o) + 1}).|a_h + b_h + a_m + b_m + a_\ell + b_\ell + b_\ell + a_\ell + b_\ell + b_$$

$$|\delta| \le ((4.u^4 + 5.u^3 + u^2).2^{-min(\alpha_o,\beta_o) + 5} + (u^2 + 2.u).2^{-min(\alpha_u + \alpha_o,\beta_u + \beta_o) + 5}).|a_h + b_h + a_m + b_m + a_\ell + b_\ell|a_h + b_h + a_m + b_m + a_\ell + b_\ell|a_h +$$

So we have $|\epsilon| \le (4.u^4 + 5.u^3 + u^2).2^{-min(\alpha_o,\beta_o)+5} + (u^2 + 2.u).2^{-min(\alpha_u + \alpha_o,\beta_u + \beta_o)+5}$

4.2 Multiplication Operators

4.2.1 Mul133

See algorithm 9

Lemma 9 (Mul133) Let a a **double** number and (b_h, b_m, b_ℓ) a **triple-Double** number, r_h , r_m and r_ℓ result of Mul133 (a, b_h, b_m, b_ℓ) for the 4 modes of rounding, considering that there is no **overflow** so: $|r_\ell| \le u.|r_m|$, $|r_m| \le (u + 2^{-\beta_o + 1} + 2^{-\beta_o - \beta_u + 1}).|r_h|$ and $|r_\ell| \le u.(u^2 + 2^{-\beta_o + 1} + 2^{-\beta_o - \beta_u + 1}).|r_h|$.

PROOF According to Mul133, we have :

 $r_m + r_\ell = Add112(t_9, t_{10})$ According to $Add112 \Rightarrow |r_\ell| \leq u.|r_m|$.

 $t_{10} = \circ(t_7 + t_8)$ According to collary 5: $t_{10} = (t_7 + t_8).(1 + \epsilon_1)$ with $|\epsilon_1| \le u$ $t_8 = \circ(t_4 + t_5)$ According to collary 5: $t_8 = (t_4 + t_5).(1 + \epsilon_2)$ with $|\epsilon_2| \le u$ $t_9 + t_7 = Add112(t_2, t_3)$ According to $Add112 \Rightarrow |t_7| \le u.|t_9|$ and $t_9 + t_7 = t_2 + t_3$.

 $t_5 = \circ(a.b_\ell)$ According to collary 5: $t_5 = (a.b_\ell).(1 + \epsilon_3)$ with $|\epsilon_3| \le u$ $t_3 + t_4 = Add112(a, b_m)$ According to $Add112 \Rightarrow |t_4| \le u.|t_3|$ and $t_3 + t_4 = a + b_m$.

 $r_h + t_2 = Mul112(a, b_h)$ According to $Mul112 \Rightarrow |t_2| \leq u.|r_h|$ and $r_h + t_2 = a.b_h$.

We search the upper bound for $|a.b_h|$ in function of $|r_h|$ to find the upper bound of $|a.b_m|$

$$|a.b_h| \le |r_h| + |t_2|$$

As $|t_2| < u.|r_h| \Rightarrow$

$$|a.b_h| \le (1+u).|r_h|$$

As $|b_m| \leq 2^{-\beta_o}.|b_h| \Rightarrow$

$$|a.b_m| \le 2^{-\beta_o}.|a.b_h|$$

 $|a.b_m| \le 2^{-\beta_o}.(1+u).|r_h|$

We have $t_3 + t_4 = a.b_m$, so:

$$|t_3 + t_4| \le 2^{-\beta_o} \cdot (1+u) \cdot |r_h|$$

As $|t_4|u.|t_3|$, so we have:

$$(1-u).|t_3| \le |t_3 + t_4|$$

 \Rightarrow

$$|t_3| \le 2^{-\beta_o} \cdot (1+u) \cdot \frac{1}{1-u} \cdot |r_h|$$

and

$$|t_4| \le 2^{-\beta_o} \cdot (1+u) \cdot \frac{u}{1-u} \cdot |r_h|$$

As
$$|b_{\ell}| \leq 2^{-\beta_u} |b_m| \Rightarrow |b_{\ell}| \leq 2^{-\beta_o - \beta_u} |b_h| \Rightarrow$$

$$|a.b_{\ell}| \le 2^{-\beta_o - \beta_u}.|a.b_h|$$

$$|a.b_{\ell}| \le 2^{-\beta_o - \beta_u} . (1+u) . |r_h|$$

As
$$t_5 = (a.b_\ell).(1 + \epsilon_3) \Rightarrow$$

$$|t_5| \le |(a.b_\ell).(1+\epsilon_3)|$$

As
$$|\epsilon_3| \leq u \Rightarrow$$

$$|t_5| \le (1+u)^2 \cdot 2^{-\beta_o - \beta_u} \cdot |r_h|$$

As $t_9 + t_7 = t_2 + t_3$ and $|t_7| \le u.|t_9|$:

$$(1-u).|t_9| \le |t_9 + t_7|$$

$$|t_9| \le \frac{1}{1-u} \cdot (|t_2| + |t_3|)$$

$$|t_9| \le \frac{1}{1-u}.(u.|r_h| + 2^{-\beta_o}.(1+u).\frac{1}{1-u}.|r_h|)$$

$$|t_9| \le \frac{1}{1-u} \cdot (u + 2^{-\beta_o} \frac{u+1}{1-u}) \cdot |r_h|$$

and so: $|t_7| \leq \frac{u}{1-u} \cdot (u + 2^{-\beta_o} \frac{u+1}{1-u}) \cdot |r_h|$ We know that $t_8 = (t_4 + t_5) \cdot (1 + \epsilon_2)$ and that $|\epsilon_2| \leq u \Rightarrow$

$$|t_8| \le (1+u).(|t_4|+|t_5|)$$

$$|t_8| \le (1+u).(2^{-\beta_o}.(1+u).\frac{u}{1-u}.|r_h| + (1+u)^2.2^{-\beta_o-\beta_u}.|r_h|)$$

$$|t_8| \le (1+u)^2 \cdot (2^{-\beta_o} \cdot \frac{u}{1-u} + (1+u) \cdot 2^{-\beta_o - \beta_u}) \cdot |r_h|$$

We have $t_{10} = (t_7 + t_8).(1 + \epsilon_1)$ and that $|\epsilon_1| \le u \Rightarrow$

$$|t_{10}| \le (1+u).(|t_7|+|t_8|)$$

$$|t_{10}| \le (1+u) \cdot \left(\frac{u}{1-u} \cdot (u+2^{-\beta_o} \frac{u+1}{1-u}) + (1+u)^2 \cdot \left(2^{-\beta_o} \cdot \frac{u}{1-u} + (1+u) \cdot 2^{-\beta_o-\beta_u}\right)\right) \cdot |r_h|$$

$$|t_{10}| \le \left(\frac{u^3 + u^2}{1 - u} + 2^{-\beta_o} - \frac{u^5 - 2u^4 + u^3 + 4u^2 + 2u}{(1 - u)^2} + 2^{-\beta_o - \beta_u} \cdot (u + 1)^4\right) \cdot |r_h|$$

And finally, $r_h + r_\ell = t_9 + t_{10}$, so:

$$|r_{m} + r_{\ell}| \leq \frac{1}{1 - u} \cdot (u + 2^{-\beta_{o}} \frac{u + 1}{1 - u}) \cdot |r_{h}|$$

$$+ (\frac{u^{3} + u^{2}}{1 - u} + 2^{-\beta_{o}} \frac{-u^{5} - 2 \cdot u^{4} + u^{3} + 4 \cdot u^{2} + 2u}{(1 - u)^{2}} + 2^{-\beta_{o} - \beta_{u}} \cdot (u + 1)^{4}) \cdot |r_{h}|$$

$$|r_{m} + r_{\ell}| \leq (\frac{u^{3} + u^{2} + u}{1 - u} + 2^{-\beta_{o}} \cdot \frac{-u^{5} - 2 \cdot u^{4} + u^{3} + 4 \cdot u^{2} + 3 \cdot u + 1}{(1 - u)^{2}} + 2^{-\beta_{o} - \beta_{u}} \cdot (1 + u)^{4}) \cdot |r_{h}|$$
As $\frac{u^{3} + u^{2} + u}{1 - u} \leq u$, $|\frac{-u^{5} - 2 \cdot u^{4} + u^{3} + 4 \cdot u^{2} + 3 \cdot u + 1}{(1 - u)^{2}}| \leq 2$ and $(1 + u)^{4} \leq 2$

$$|r_{m} + r_{\ell}| \leq (u + 2^{-\beta_{o} + 1} + 2^{-\beta_{o} - \beta_{u} + 1}) \cdot |r_{h}|$$
As $|r_{\ell}| \leq u |r_{m}| \Rightarrow$

As
$$|r_{\ell}| \le u.|r_{m}| \Rightarrow$$

$$(1-u).|r_{m}| \le |r_{m}+r_{\ell}|$$

$$(1-u).|r_{m}| \le (u+2^{-\beta_{o}+1}+2^{-\beta_{o}-\beta_{u}+1}).|r_{h}|$$

$$|r_{m}| \le |\frac{1}{1-u}|.(u+2^{-\beta_{o}+1}+2^{-\beta_{o}-\beta_{u}+1}).|r_{h}|$$

$$|r_{m}| \le (u+2^{-\beta_{o}+1}+2^{-\beta_{o}-\beta_{u}+1}).|r_{h}|$$

and

$$|r_{\ell}| \le u.(u^2 + 2^{-\beta_o + 1} + 2^{-\beta_o - \beta_u + 1}).|r_h|$$

Theorem 10 (Relative error algorithm Mul133) Let a a double number and (b_h, b_m, b_ℓ) a triple-double number are the arguments of the function Mul133.

So:

$$r_h + r_m + r_\ell = (a.(b_h + b_m + b_\ell)).(1 + \epsilon)$$
with $|\epsilon| \le \frac{u^4 + 4.u^3.2^{-\beta_o} + 4.u.2^{-\beta_o - \beta_u}}{1 - (2^{-\beta_o} + 2^{-\beta_o - \beta_u})} \le 2.u^4 + u^3.2^{-\beta_o + 3} + u.2^{-\beta_o - \beta_u + 3}.$

PROOF According to the calculation technique of the theorem 5.2 prof's([12]). According to Add333 and thanks to Add112 and Mul112, we have these results:

$$|t_2| \le u.|r_h|$$

 $|t_4| \le u.|t_3|$
 $|t_7| \le u.|t_9|$
 $|r_\ell| \le u.|r_m|$

We search for all $|t_i|$ with $2 \le i \le 10$, their upper bound in function of $|a.b_h|$.

As $r_h + t_2 = a.b_h$ and $|t_2| \le u.|r_h|$, we have $(1-u).|r_h| \le |r_h + t_2| \Rightarrow$

$$(1-u).|r_h| < |a.b_h|$$

$$|r_h| \le \left| \frac{1}{1 - u} . |a.b_h| \right|$$

 \Rightarrow

$$|t_2| \le |\frac{u}{1-u}.|a.b_h|$$

As $|b_m| \le 2^{-\beta_o} . |b_h| \Rightarrow |a.b_m| \le 2^{-\beta_o} . |a.b_h| \Rightarrow$

$$t_3 + t_4 = a.b_m$$

As $|t_4| \le u.|t_3| \Rightarrow (1-u).|t_3| \le |t_3+t_4|$

$$(1-u).|t_3| \le |a.b_m|$$

 $|t_3| \leq \frac{1}{1-a}.2^{-\beta_o}.|a.b_h|$

 \Rightarrow

$$|t_4| \leq \frac{u}{1-u}.2^{-\beta_o}.|a.b_h|$$

We have $t_5 = \circ(a.b_\ell)$ After the collary $t_5 = a.b_\ell.(1 + \epsilon_3)$ with $|\epsilon_3| \le u$. As $|b_\ell| \le 2^{-\beta_u}.|b_m| \le 2^{-\beta_o-\beta_u}.|b_h| \Rightarrow |a.b_\ell| \le 2^{-\beta_o-\beta_u}.|a.b_h|$

$$t_5 = a.b_{\ell}.(1 + \epsilon_3)$$

$$|t_5| < |a.b_{\ell}|.(1+u)$$

$$|t_5| \le (1+u).2^{-\beta_o - \beta_u}.|a.b_h|$$

After calculation, we have:

$$|t_{9}| \leq (u + 2^{-\beta_{o}}).|\frac{1}{1 - u}.|a.b_{h}|$$

$$|t_{7}| \leq (u + 2^{-\beta_{o}}).|\frac{u}{1 - u}.|a.b_{h}|$$

$$|t_{8}| \leq (\frac{u^{3} + u^{2}}{1 - u}.2^{-\beta_{o}} + (1 + u)^{2}.2^{-\beta_{o} - \beta_{u}}).|a.b_{h}|$$

$$|t_{10}| \le \frac{1}{u-1} (u^3 + u^2 + (u^4 + 3.u^3 + 2.u^2).2^{-\beta_o} + (u^4 + 2.u^3 - 2.u - 1).2^{-\beta_o - \beta_u}).|a.b_h|$$

We know that $r_m + r_\ell = t_9 + t_{10}$, we start by calculating t_{10} :

$$t_{10} = (t_7 + t_8)(1 + \epsilon_1)$$

with $|\epsilon_1| \leq u$, then we calculate t_8 :

$$t_8 = (t_4 + t_5).(1 + \epsilon_2)$$

with $|\epsilon_2| \leq u$, then we calculate t_5 :

$$t_5 = a.b_{\ell}.(1 + \epsilon_3)$$

with $|\epsilon_3| \leq u$

$$t_8 = (t_4 + a.b_{\ell}.(1 + \epsilon_3)).(1 + \epsilon_2)$$

$$t_{10} = (t_7 + (t_4 + a.b_{\ell}.(1 + \epsilon_3)).(1 + \epsilon_2))(1 + \epsilon_1)$$

After the calculation, we have:

$$t_{10} = t_7 + t_4 + a.b_{\ell} + \delta$$

with $\delta = t_7 \cdot \epsilon_1 + (t_4 + a \cdot b_\ell) \cdot (\epsilon_3 + \epsilon_2 + \epsilon_3 \cdot \epsilon_2 + \epsilon_1 + \epsilon_3 \cdot \epsilon_1 + \epsilon_2 \cdot \epsilon_1 + \epsilon_3 \cdot \epsilon_2 \cdot \epsilon_1)$ As $r_h + r_m + r_\ell = r_h + t_9 + t_{10}$, we have:

$$r_b + r_m + r_\ell = r_b + t_0 + t_7 + t_4 + a.b_\ell + \delta$$

As $t_9 + t_7 = t_2 + t_3 \Rightarrow$

$$r_h + r_m + r_\ell = r_h + t_2 + t_3 + t_4 + a.b_\ell + \delta$$

As $r_h + t_3 = a.b_h$ and $t_3 + t_4 = a.b_m \Rightarrow$

$$r_h + r_m + r_\ell = a.b_h + a.b_m + a.b_\ell + \delta$$

$$r_h + r_m + r_\ell = a.(b_h + b_m + b_\ell) + \delta$$

We seek the upper bound of $|\delta|$ in function of $|a.b_h|$, So:

$$|\delta| \le |t_7.\epsilon_1| + (|t_4| + |a.b_\ell|).|\epsilon_3 + \epsilon_2 + \epsilon_3.\epsilon_2 + \epsilon_1 + \epsilon_3.\epsilon_1 + \epsilon_2.\epsilon_1 + \epsilon_3.\epsilon_2.\epsilon_1|$$

As $|\epsilon_i| < u$ with $1 < i < 3 \Rightarrow$

$$|\delta| \le u.|t_7| + (3.u + 3.u^2 + u^3).(|t_4| + |a.b_\ell|)$$

$$|\delta| \le u.(u+2^{-\beta_o}).|\frac{u^2}{u-1}.|a.b_h| + (3.u+3.u^2+u^3).(\frac{u^2}{u-1}.2^{-\beta_o}.|a.b_h| + 2^{-\beta_o-\beta_u}.|a.b_h|)$$

$$|\delta| \le \left(\frac{u^4}{u-1} + \frac{u^5 + 3u^4 + 4u^3}{u-1} \cdot 2^{-\beta_o} + \left(u^3 + 3u^2 + 3u\right) \cdot 2^{-\beta_o - \beta_u}\right) \cdot |a \cdot b_h|$$

We seek the upper bound of $|a.b_h|$ in function of $|a.(b_h + b_m + b_\ell)|$.

$$|a.b_m + a.b_{\ell}| \le |a.b_m| + |a.b_{\ell}|$$

$$|a.b_m + a.b_{\ell}| \le 2^{-\beta_o}.|a.b_h| + 2^{-\beta_o - \beta_u}.|a.b_h|$$

$$|a.b_m + a.b_{\ell}| \le (2^{-\beta_o} + 2^{-\beta_o - \beta_u}).|a.b_h|$$

So we have:

$$|a.b_h + a.b_m + a.b_\ell| \ge (1 - (2^{-\beta_o} + 2^{-\beta_o - \beta_u})) \cdot |a.b_h|$$

$$\frac{1}{1 - (2^{-\beta_o} + 2^{-\beta_o - \beta_u})} |a.b_h + a.b_m + a.b_\ell| \ge |a.b_h|$$

 \Rightarrow

$$|\delta| \le \left(\frac{u^4}{u-1} + \frac{u^5 + 3.u^4 + 4.u^3}{u-1} \cdot 2^{-\beta_o} + \left(u^3 + 3.u^2 + 3.u\right) \cdot 2^{-\beta_o - \beta_u}\right) \cdot \frac{1}{1 - \left(2^{-\beta_o} + 2^{-\beta_o - \beta_u}\right)} |a.b_h + a.b.$$

As
$$\frac{u^4}{u-1} \le u^4$$
, $\frac{u^5 + 3.u^4 + 4.u^3}{u-1} \le 4.u^3$ and $u^3 + 3.u^2 + 3.u \le 4.u \Rightarrow$

$$|\delta| \le \frac{u^4 + 4.u^3 \cdot 2^{-\beta_o} + 4.u \cdot 2^{-\beta_o - \beta_u}}{1 - (2^{-\beta_o} + 2^{-\beta_o - \beta_u})} . |a.b_h + a.b_m + a.b_\ell|$$

As
$$1 - (2^{-\beta_o} + 2^{-\beta_o - \beta_u}) \ge \frac{1}{2}$$
 with $\beta_o \ge 2$ and $\beta_u \ge 2$

$$|\delta| \le (2.u^4 + u^3.2^{-\beta_o + 3} + u.2^{-\beta_o - \beta_u + 3}).|a.b_h + a.b_m + a.b_\ell|$$

So:

$$|\epsilon| \le 2.u^4 + u^3 \cdot 2^{-\beta_o + 3} + u \cdot 2^{-\beta_o - \beta_u + 3}$$

4.2.2 Mul233

See algorithm 10

Lemma 10 (Mul233) Let (a_h, a_ℓ) a double-double number and (b_h, b_m, b_ℓ) a triple-Double number, r_h , r_m and r_ℓ result of Mul233 $(a_h, a_\ell, b_h, b_m, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so: $|r_\ell| \leq 4.u.|r_m|$, $|r_m| \leq u.(2.u + 2^{-\beta_o} + 2^{-\beta_o-\beta_u}).|r_h|$ and $|r_\ell| \leq u.(2.u + 2^{-\beta_o} + 2^{-\beta_o-\beta_u}).|r_h|$.

Proof According to Mul233, we have:

• $r_h + t_1 = Mul112(a_h, b_h)$ according to $Mul112 \Rightarrow r_h + t_1 = a_h.b_h$ and $|t_1| \leq u.|r_h| \Rightarrow$

$$|a_h.b_h| \le (1+u).|r_h|$$

 $|a_h.b_h| \le \frac{1}{u+1}.|r_h|$

• $t_2 + t_3 = Mul112(a_h, b_m)$ according to $Mul112 \Rightarrow t_2 + t_3 = a_h.b_m$ and $|t_3| \leq u.|t_2| \Rightarrow$

$$|t_2 + t_3| \le |a_h \cdot b_m|$$
As $|b_m| \le 2^{-\beta_o} \cdot |b_h| \Rightarrow |a_h \cdot b_m| \le 2^{-\beta_o} \cdot |a_h \cdot b_h| \Rightarrow$
 $|t_2 + t_3| \le 2^{-\beta_o} \cdot |a_h \cdot b_h|$

As $|a_h.b_h| \leq \frac{1}{u+1}.|r_h|$

$$|t_2 + t_3| \le 2^{-\beta_o} \cdot \frac{1}{u+1} \cdot |r_h|$$

As
$$|t_3| \le u.|t_2| \Rightarrow$$

$$(1-u).|t_2| \le |t_2 + t_3|$$

$$(1-u).|t_2| \le 2^{-\beta_o}.\frac{1}{u+1}.|r_h|$$

$$|t_2| \le 2^{-\beta_o}.\frac{1}{1-u^2}.|r_h|$$

$$\Rightarrow$$

$$|t_3| \le 2^{-\beta_o}.\frac{u}{1-u^2}.|r_h|$$

So not to write the repetitive calculations, I give the results for $|t_i|$ with $4 \le i \le 9$:

$$|t_4| \le 2^{-\beta_o - \beta_u} \cdot \frac{1}{1 - u^2} \cdot |r_h|$$

$$|t_5| \le 2^{-\beta_o - \beta_u} \cdot \frac{u}{1 - u^2} \cdot |r_h|$$

$$|t_6| \le \frac{u}{1 - u^2} \cdot |r_h|$$

$$|t_7| \le \frac{u^2}{1 - u^2} \cdot |r_h|$$

$$|t_8| \le 2^{-\beta_o} \cdot \frac{u}{1 - u^2} \cdot |r_h|$$

$$|t_9| \le 2^{-\beta_o} \cdot \frac{u^2}{1 - u^2} \cdot |r_h|$$

• $t_{10} = \circ(a_{\ell}.b_{\ell})$ According to collary 5: $t_{10} = a_{\ell}.b_{\ell}(1 + \epsilon_5)$ with $|\epsilon_5| \leq u$ but $|a_{\ell}| \leq u.|a_h|$ and $|b_{\ell}| \leq 2^{-\beta_o - \beta_u}.|b_h| \Rightarrow$

$$|t_{10}| \le u.2^{-\beta_o - \beta_u}.|a_h.b_h|.(1+u)$$

$$|t_{10}| \le u.(1+u).2^{-\beta_o - \beta_u}.\frac{1}{u+1}.|r_h|$$

$$|t_{10}| \le u.2^{-\beta_o - \beta_u}.|r_h|$$

• $t_{11} + t_{12} = Add222(t_2, t_3, t_4, t_5)$ according to $Add222 \Rightarrow t_{11} + t_{12} = (t_2 + t_3 + t_4 + t_5).(1 + \epsilon_4)$ with $|\epsilon_4| \le 6.u^2$ and $|t_{12}| \le 6.u.|t_{11}| \Rightarrow$

$$|t_{11} + t_{12}| \ge (1 - 6.u).|t_{11}|$$

$$(1 - 6.u).|t_{11}| \le (|t_2| + |t_3| + |t_4| + |t_5|).(1 + 6.u^2)$$
$$|t_{11}| \le \frac{1 + 6.u^2}{1 - 6.u}(|t_2| + |t_3| + |t_4| + |t_5|)$$

$$|t_{11}| \leq \frac{1+6.u^2}{1-6.u} \cdot (2^{-\beta_o} \cdot \frac{1}{1-u^2} \cdot |r_h| + 2^{-\beta_o} \cdot \frac{u}{1-u^2} \cdot |r_h| + 2^{-\beta_o-\beta_u} \cdot \frac{1}{1-u^2} \cdot |r_h|$$

$$+2 \cdot 2^{-\beta_o-\beta_u} \cdot \frac{u}{1-u^2} \cdot |r_h|)$$

$$|t_{11}| \leq \frac{(1+6.u^2) \cdot (u+1)}{(1-6.u) \cdot (1-u^2)} \cdot (2^{-\beta_o} + 2^{-\beta_o-\beta_u}) \cdot |r_h|$$

$$|t_{12}| \leq \frac{(1+6.u^2) \cdot (u^2+u)}{(1-6.u) \cdot (1-u^2)} \cdot (2^{-\beta_o} + 2^{-\beta_o-\beta_u}) \cdot |r_h|$$

• $t_{13} + t_{14} = Add222(t_6, t_7, t_8, t_9)$ according to $Add222 \Rightarrow t_{13} + t_{14} = (t_6 + t_7 + t_8 + t_9).(1 + \epsilon_3)$ with $|\epsilon_3| \leq 6.u^2$ and $|t_{14}| \leq 6.u.|t_{13}|$ We do the same calculations as before.

$$\Rightarrow |t_{13}| \leq \frac{1}{1 - 6.u} \cdot (|t_6| + |t_7| + |t_8| + |t_9|) \cdot (1 + 6.u^2)$$

$$|t_{13}| \leq \frac{1 + 6.u^2}{1 - 6.u} \cdot (\frac{u}{1 - u^2} \cdot |r_h| + \frac{u^2}{1 - u^2} \cdot |r_h| + 2^{-\beta_o} \cdot \frac{u}{1 - u^2} \cdot |r_h| + 2^{-\beta_o} \cdot \frac{u^2}{1 - u^2} \cdot |r_h|)$$

$$|t_{13}| \leq \frac{(1 + 6.u^2) \cdot (u^2 + u) \cdot (u + 1)^2}{(1 - 6.u) \cdot (1 - u^2)} \cdot (1 + 2^{-\beta_o}) \cdot |r_h|.$$

$$|t_{14}| \le \frac{(1+6.u^2).(u^3+u^2).(u+1)^2}{(1-6.u).(1-u^2)}.(1+2^{-\beta_o}).|r_h|.$$

we do the same calculation for t_{15} and t_{16} , after the calculations then the simplifications, we have:

$$|t_{15}| \leq \frac{(1+6.u^2).(u+1)^2}{(1-6.u).(1-u^2)}.(u.(u+1)^2 + (u^3+2.u^2+u+1).2^{-\beta_o} + 2^{-\beta_o-\beta_u}).|r_h|$$

$$|t_{16}| \leq \frac{u.(1+6.u^2).(u+1)^2}{(1-6.u).(1-u^2)}.(u.(u+1)^2 + (u^3+2.u^2+u+1).2^{-\beta_o} + 2^{-\beta_o-\beta_u}).|r_h|$$

• $t_{17} + t_{18} = Add112(t_1, t_{10})$ according to $Add112 \Rightarrow t_{17} + t_{18} = t_1 + t_{10}$ and $|t_{18}| \leq u.|t_{17}|$, as done before \Rightarrow

$$|t_{17}| \le \left| \frac{1}{1-u} \right| . (|t_1| + |t_{10}|)$$

$$|t_{17}| \le \left| \frac{1}{1-u} \right| . (u.|r_h| + u.2^{-\beta_o - \beta_u}.|r_h|)$$

$$|t_{17}| \le \left| \frac{u}{1-u} \right| . (|r_h| + 2^{-\beta_o - \beta_u}.|r_h|)$$

$$|t_{18}| \le \left| \frac{u^2}{1-u} \right| . (|r_h| + 2^{-\beta_o - \beta_u}.|r_h|)$$

• $r_m + r_\ell = Add222(t_{17}, t_{18}, t_{15}, t_{16})$ according to $Add222 \Rightarrow r_m + r_\ell = (t_{17} + t_{18} + t_{15} + t_{16})(1 + \epsilon_1)$ with $|\epsilon_1| \leq 6.u^2$ and $|r_\ell| \leq 6.u.|r_m|$, after calculation:

$$|r_m| \le \frac{1}{1-6u} \cdot (|t_{17}| + |t_{18}| + |t_{15}| + |t_{16}|) \cdot (1+6u^2)$$

after we calculate:

$$|r_m| \le (2.u + 2^{-\beta_o} + 2^{-\beta_o - \beta_u}).|r_h|$$

and

$$|r_{\ell}| \le u.(2.u + 2^{-\beta_o} + 2^{-\beta_o - \beta_u}).|r_h|$$

Theorem 11 (Relative error algorithm Mul233) Let (a_h, a_ℓ) a double-double number and (b_h, b_m, b_ℓ) a triple-double number are the arguments of the function Mul233. So:

$$r_h + r_m + r_\ell = (a_h + a_\ell).(b_h + b_m + b_\ell)).(1 + \epsilon)$$

$$with |\epsilon| \le \frac{24.u^3 + 19.u^2.2^{-\beta_o} + 20.u^2.2^{-\beta_o - \beta_u}}{1 - (u + (1 + u).2^{-\beta_o} + (1 + u).2^{-\beta_o - \beta_u}} \le 48.u^3 + 19.u^2.2^{-\beta_o + 1} + 5.u^2.2^{-\beta_o - \beta_u + 3}$$

PROOF (10PT) According to the conditions of Mul233, we begin to calculate $r_m + r_\ell$:

 $r_m + r_\ell = Add222(t_{17}, t_{18}, t_{15}, t_{16})$ According to $Add222 \Rightarrow r_m + r_\ell = (t_{17} + t_{18} + t_{15} + t_{16}).(1 + \epsilon_1)$ with $|\epsilon_1| \leq 6.u^2$ As $t_{17} + t_{18} = t_1 + t_{10}$ and $t_{15} + t_{16} = (t_{11} + t_{12} + t_{13} + t_{14}).(1 + \epsilon_2)$ with $|\epsilon_2| \leq 6.u^2$

$$r_m + r_\ell = (t_1 + t_{10} + (t_{11} + t_{12} + t_{13} + t_{14}).(1 + \epsilon_2)).(1 + \epsilon_1)$$

As $t_{13} + t_{14} = (t_6 + t_7 + t_8 + t_9)(1 + \epsilon_3)$ and $t_{11} + t_{12} = (t_2 + t_3 + t_4 + t_5)(1 + \epsilon_4)$ with $|\epsilon_2| \le 6 \cdot u^2$ and $|\epsilon_3| \le 6 \cdot u^2 \Rightarrow$

$$r_m + r_\ell = (t_1 + t_{10} + ((t_2 + t_3 + t_4 + t_5)(1 + \epsilon_4) + (t_6 + t_7 + t_8 + t_9)(1 + \epsilon_3)) \cdot (1 + \epsilon_2)) \cdot (1 + \epsilon_1)$$

As
$$t_2 + t_3 = a_h \cdot b_m$$
, $t_4 + t_5 = a_h \cdot b_\ell$, $t_6 + t_7 = a_\ell \cdot b_h$ and $t_8 + t_9 = a_\ell \cdot b_m \Rightarrow$

$$r_m + r_\ell = (t_1 + t_{10} + ((a_h \cdot b_m + a_h \cdot b_\ell)(1 + \epsilon_4) + (a_\ell \cdot b_h + a_\ell \cdot b_m)(1 + \epsilon_3)) \cdot (1 + \epsilon_2)) \cdot (1 + \epsilon_1)$$

As
$$t_{10} = a_{\ell}.b_{\ell}.(1 + \epsilon_5)$$
 with $|\epsilon_5| \le u \Rightarrow$

$$r_m + r_\ell = (t_1 + a_\ell . b_\ell . (1 + \epsilon_5) + ((a_h . b_m + a_h . b_\ell) (1 + \epsilon_4) + (a_\ell . b_h + a_\ell . b_m) (1 + \epsilon_3)) . (1 + \epsilon_2)) . (1 + \epsilon_1)$$

As
$$r_h + t_1 = a_h.b_h \Rightarrow$$

$$r_h + r_m + r_\ell = (a_h + a_\ell) \cdot (b_h + b_m + b_\ell) + \delta$$

with $\delta = t_1 \cdot \epsilon_1 + a_\ell \cdot b_\ell \cdot (\epsilon_1 + \epsilon_5 + \epsilon_1 \cdot \epsilon_5) + (a_h \cdot b_m + a_h \cdot b_\ell) \cdot (\epsilon_1 + \epsilon_2 \epsilon_4 + \epsilon_1 \cdot \epsilon_2 + \epsilon_1 \cdot \epsilon_4 + \epsilon_2 \cdot \epsilon_4 + \epsilon_1 \cdot \epsilon_2 \cdot \epsilon_4) + (a_\ell \cdot b_h + a_\ell \cdot b_m) \cdot (\epsilon_1 + \epsilon_3 \epsilon_4 + \epsilon_1 \cdot \epsilon_3 + \epsilon_1 \cdot \epsilon_4 + \epsilon_3 \cdot \epsilon_4 + \epsilon_1 \cdot \epsilon_3 \cdot \epsilon_4).$ We seek the upper bound of $|\delta|$ in function of $|a_h \cdot b_h|$:

$$|\delta| \leq |t_1.\epsilon_1| + |a_\ell.b_\ell.(\epsilon_1 + \epsilon_5 + \epsilon_1.\epsilon_5)| + (|a_h.b_m| + |a_h.b_\ell|).|\epsilon_1 + \epsilon_2\epsilon_4 + \epsilon_1.\epsilon_2 + \epsilon_1.\epsilon_4 + \epsilon_2.\epsilon_4 + \epsilon_1.\epsilon_2.\epsilon_4|$$

$$+(|a_{\ell}.b_h|+|a_{\ell}.b_m|).|\epsilon_1+\epsilon_3\epsilon_4+\epsilon_1.\epsilon_3+\epsilon_1.\epsilon_4+\epsilon_3.\epsilon_4+\epsilon_1.\epsilon_3.\epsilon_4|$$

$$|\delta| \le 6.u^2.|t_1| + (6.u^3 + 6.u^2 + u).|a_\ell.b_\ell.| + (216.u^6 + 108.u^4 + 18.u^2).(|a_h.b_m| + |a_h.b_\ell| + |a_\ell.b_h| + |a_\ell.$$

We seek the upper bound of $|a_h.b_m|$, $|a_h.b_\ell|$, $|a_\ell.b_h|$, $|a_\ell.b_m|$ and $|a_\ell.b_\ell|$ in function of $|a_h.b_h|$.

We have:

$$|a_h.b_m| \le 2^{-\beta_o}.|a_h.b_h|$$

$$|a_h.b_\ell| \le 2^{-\beta_o-\beta_u}.|a_h.b_h|$$

$$|a_\ell.b_h| \le u.|a_h.b_h|$$

$$|a_\ell.b_m| \le u.2^{-\beta_o}.|a_h.b_h|$$

$$|a_\ell.b_\ell| \le u.2^{-\beta_o-\beta_u}.|a_h.b_h|$$

We search the upper bound of $(|a_h.b_m| + |a_h.b_\ell| + |a_\ell.b_h| + |a_\ell.b_m|)$ in function of $|a_h.b_h|$

$$|a_h.b_m| + |a_h.b_\ell| + |a_\ell.b_h| + |a_\ell.b_m| \le 2^{-\beta_o}.|a_h.b_h| + 2^{-\beta_o-\beta_u}.|a_h.b_h| + u.|a_h.b_h| + u.2^{-\beta_o}.|a_h.b_h|$$

$$|a_h.b_m| + |a_h.b_\ell| + |a_\ell.b_h| + |a_\ell.b_m| \le (u + (1+u).2^{-\beta_o} + 2^{-\beta_o-\beta_u}).|a_h.b_h|$$

as $|a_{\ell}.b_{\ell}| \le u.2^{-\beta_o - \beta_u}.|a_h.b_h|$

According to Mul112, we have :

$$|t_1| \le u.|r_h| \ and \ r_h + t_1 = a_h.b_h$$

We begin to calculate the lower bound of $|r_h + t_1|$ compared to $|r_h|$. As $|t_1| \leq u.|r_h|$, we have:

$$|r_h + t_1| \ge (1 - u).|r_h|$$

 \Rightarrow

$$(1-u).|r_h| \le |a_h.b_h|$$

$$|r_h| \leq \frac{1}{1-a}.|a_h.b_h|$$

 \Rightarrow

$$|t_1| \leq \frac{u}{1-u}.|a_h.b_h|$$

$$|\delta| \le \left(\frac{6u^3}{1-u} \cdot + (6.u^3 + 6.u^2 + u) \cdot u \cdot 2^{-\beta_o - \beta_u} + (216.u^6 + 108.u^4 + 18.u^2) \cdot (u + (1+u) \cdot 2^{-\beta_o} + 2^{-\beta_o - \beta_u})\right) \cdot |a_h \cdot b_h|$$

we simplify with the upper bounds of each coefficient.

$$|\delta| \le (24.u^3 + 19.u^2.2^{-\beta_o} + 20.u^2.2^{-\beta_o-\beta_u}).|a_h.b_h|$$

Now, we seek the lower bound of $|a_h.b_h|$ in function of $|(a_h+a_\ell).(b_h+b_m+b_\ell)|$, as $|a_h.b_m| + |a_h.b_\ell| + |a_\ell.b_h| + |a_\ell.b_m| \le (u + (1+u).2^{-\beta_o} + 2^{-\beta_o-\beta_u}).|a_h.b_h|$ and $|a_\ell.b_\ell| \le u.2^{-\beta_o-\beta_u}.|a_h.b_h| \Rightarrow$

$$|a_h.b_m+a_h.b_\ell+a_\ell.b_h+a_\ell.b_m+a_\ell.b_\ell| \le |a_h.b_m|+|a_h.b_\ell|+|a_\ell.b_h|+|a_\ell.b_m|+|a_\ell.b_\ell|$$

$$|a_h.b_m + a_h.b_\ell + a_\ell.b_h + a_\ell.b_m + a_\ell.b_\ell| \le (u + (1+u).2^{-\beta_o} + (1+u).2^{-\beta_o-\beta_u}).|a_h.b_h|$$

So we have:

$$|(a_h + a_\ell) \cdot (b_h + b_m + b_\ell)| \ge (1 - (u + (1+u) \cdot 2^{-\beta_o} + (1+u) \cdot 2^{-\beta_o - \beta_u})) \cdot |a_h \cdot b_h|$$

$$|a_h.b_h| \le \frac{1}{1 - (u + (1+u).2^{-\beta_o} + (1+u).2^{-\beta_o - \beta_u})}.|(a_h + a_\ell).(b_h + b_m + b_\ell)|$$

at the end:

$$|\delta| \leq \frac{24 \cdot u^3 + 19 \cdot u^2 \cdot 2^{-\beta_o} + 20 \cdot u^2 \cdot 2^{-\beta_o - \beta_u}}{1 - (u + (1 + u) \cdot 2^{-\beta_o} + (1 + u) \cdot 2^{-\beta_o - \beta_u})} . |(a_h + a_\ell) \cdot (b_h + b_m + b_\ell)|$$

$$\Rightarrow |\epsilon| \leq \frac{24 \cdot u^3 + 19 \cdot u^2 \cdot 2^{-\beta_o + 20 \cdot u^2} \cdot 2^{-\beta_o - \beta_u}}{1 - (u + (1 + u) \cdot 2^{-\beta_o - \beta_u})}$$
As $\beta_o \geq 2$ and $\beta_u \geq 1 \Rightarrow 1 - (u + (1 + u) \cdot 2^{-\beta_o} + (1 + u) \cdot 2^{-\beta_o - \beta_u}) \geq \frac{1}{2}$

$$\Rightarrow |\epsilon| \leq 48 \cdot u^3 + 19 \cdot u^2 \cdot 2^{-\beta_o + 1} + 5 \cdot u^2 \cdot 2^{-\beta_o - \beta_u + 3}$$

4.2.3 Mul333

See algorithm 11

Lemma 11 (Mul333) Let (a_h, a_m, a_ℓ) and (b_h, b_m, b_ℓ) are **triple-Double** number, r_h , r_m and r_ℓ result of Mul333 $(a_h, a_m, a_\ell, b_h, b_m, b_\ell)$ for the 4 modes of rounding, considering that there is no **overflow** so: $|r_\ell| \leq 6.u.|r_m|$, $|r_m| \leq .(u + 2^{-\alpha_o} + 2^{-\beta_o} + 2^{-\alpha_o - \beta_o + 1}).|r_h|$ and $|r_\ell| \leq u.(u + 2^{-\alpha_o} + 2^{-\beta_o} + 2^{-\alpha_o - \beta_o + 1}).|r_h|$

PROOF According to Mul333, we seek for each $|t_i|$ with $1 \le i \le 22$ their upper bound in function of $|r_h|$

We begin to search the upper bound of $|a_h.b_m|$, $|a_h.b_\ell|$, $|a_m.b_h|$, $|a_m.b_m|$, $|a_m.b_\ell|$, $|a_\ell.b_m|$ and $|a_\ell.b_\ell|$ in function of $|a_h.b_h|$. Thanks to the conditions of Add333, we have:

$$|a_{h}.b_{m}| \leq 2^{-\beta_{o}}.|a_{h}.b_{h}|$$

$$|a_{h}.b_{\ell}| \leq 2^{-\beta_{o}-\beta_{u}}.|a_{h}.b_{h}|$$

$$|a_{m}.b_{h}| \leq 2^{-\alpha_{o}}.|a_{h}.b_{h}|$$

$$|a_{m}.b_{m}| \leq 2^{-\alpha_{o}-\beta_{o}}.|a_{h}.b_{h}|$$

$$|a_{m}.b_{\ell}| \leq 2^{-\alpha_{o}-\beta_{o}-\beta_{u}}.|a_{h}.b_{h}|$$

$$|a_{\ell}.b_{h}| \leq 2^{-\alpha_{o}-\alpha_{u}}.|a_{h}.b_{h}|$$

$$|a_{\ell}.b_{m}| \leq 2^{-\alpha_{o}-\alpha_{u}-\beta_{o}}.|a_{h}.b_{h}|$$

$$|a_{\ell}.b_{\ell}| < 2^{-\alpha_{o}-\alpha_{u}-\beta_{o}-\beta_{u}}.|a_{h}.b_{h}|$$

$$|a_{\ell}.b_{\ell}| < 2^{-\alpha_{o}-\alpha_{u}-\beta_{o}-\beta_{u}}.|a_{h}.b_{h}|$$

We have:

• $r_h + t_1 = Mul112(a_h, b_h)$ according to $Mul112 \Rightarrow r_h + t_1 = a_h.b_h$ and $|t_1| \leq u.|r_h| \Rightarrow$

$$|a_h.b_h| \le (1+u).|r_h|$$

• $t_2 + t_3 = Mul112(a_h, b_m)$ according to $Mul112 \Rightarrow t_2 + t_3 = a_h.b_m$ and $|t_3| \leq u.|t_2| \Rightarrow$

$$|t_{2} + t_{3}| \ge (1 - u).|t_{2}|$$

$$(1 - u).|t_{2}| \le |a_{h}.b_{m}|$$

$$|t_{2}| \le \frac{1}{1 - u}.|a_{h}.b_{m}|$$

$$|t_{2}| \le \frac{1}{1 - u}.2^{-\beta_{o}}.|a_{h}.b_{h}|$$

$$|t_{2}| \le \frac{1}{1 - u}.2^{-\beta_{o}}.(1 + u).|r_{h}|$$

$$|t_{2}| \le \frac{1 + u}{1 - u}.2^{-\beta_{o}}.|r_{h}|$$

$$|t_{3}| \le \frac{u^{2} + u}{1 - u}.2^{-\beta_{o}}.|r_{h}|$$

 \Rightarrow

For $|t_i|$ with $4 \le i \le 7$, we don't repeat the calculations:

$$|t_4| \le \frac{1+u}{1-u} \cdot 2^{-\alpha_o} \cdot |r_h|$$

$$|t_5| \le \frac{u^2+u}{1-u} \cdot 2^{-\alpha_o} \cdot |r_h|$$

$$|t_6| \le \frac{1+u}{1-u} \cdot 2^{-\alpha_o-\beta_o} \cdot |r_h|$$

$$|t_7| \le \frac{u^2+u}{1-u} \cdot 2^{-\alpha_o-\beta_o} \cdot |r_h|$$

• $t_8 = \circ(a_h.b_\ell)$ according to the collary $5 \Rightarrow t_8 = a_h.b_\ell.(1 + \epsilon_{11})$ with $|\epsilon_{11}| \leq u \Rightarrow$

$$|t_8| \le |a_h.b_\ell|.(1+u)$$

$$|t_8| \le 2^{-\beta_o-\beta_u}.|a_h.b_h|.(1+u)$$

$$|t_8| \le 2^{-\beta_o-\beta_u}.(1+u).|r_h|.(1+u)$$

$$|t_8| \le (1+u)^2.2^{-\beta_o-\beta_u}.|r_h|$$

We do the same operations for $|t_i|$ with $9 \le i \le 11$:

$$|t_9| \le (1+u)^2 \cdot 2^{-\alpha_o - \alpha_u} \cdot |r_h|$$

$$|t_{10}| \le (1+u)^2 \cdot 2^{-\alpha_o - \beta_o - \beta_u} \cdot |r_h|$$

$$|t_{11}| \le (1+u)^2 \cdot 2^{-\alpha_o - \alpha_u - \beta_o} \cdot |r_h|$$

• $t_{12} = \circ(t_8 + t_9)$ according to collary 5: $t_{12} = (t_8 + t_9).(1 + \epsilon_7)$ with $|\epsilon_7| \le u \Rightarrow$

$$|t_{12}| \le |t_8 + t_9|.(1+u)$$

$$|t_{12}| \le ((1+u)^2.2^{-\beta_o - \beta_u} + (1+u)^2.2^{-\alpha_o - \alpha_u}).|r_h|.(1+u)$$

$$|t_{12}| \le ((1+u)^3.2^{-\min(\beta_o + \beta_u, \alpha_o + \alpha_u)} + .|r_h|$$

similarly:

$$|t_{13}| \le ((1+u)^3 \cdot 2^{-\min(\alpha_o + \beta_o + \beta_u, \alpha_o + \alpha_u + \beta_o)} |r_h|$$

• $t_{14} + t_{15} = Add112(t_1, t_6)$ according to $Add112 \Rightarrow t_{14} + t_{15} = t_1 + t_6$ and $|t_{15}| \le u.|t_{14}| \Rightarrow$

$$|t_{14}| \le \frac{1}{1-u} \cdot |t_1 + t_6|$$

$$|t_{14}| \le \frac{1}{1-u} \cdot (u + \frac{1+u}{1-u} \cdot 2^{-\alpha_o - \beta_o}) \cdot |r_h|$$

$$|t_{14}| \le \frac{1}{(1-u)^2} \cdot (u - u^2 + (1+u) \cdot 2^{-\alpha_o - \beta_o}) \cdot |r_h|$$

$$\Rightarrow |t_{15}| \le \frac{u}{(1-u)^2} \cdot (u - u^2 + (1+u) \cdot 2^{-\alpha_o - \beta_o}) \cdot |r_h|$$

• We calculate $|t_i|$ with $16 \le i \le 18$ as we did previously with the calculation of $|t_{12}|$

$$|t_{16}| \leq \frac{1+u}{(1-u)^2} \cdot (u^2 - u^3 + (-u^3 + u) \cdot 2^{-\alpha_o - \beta_o}) \cdot |r_h|$$

$$|t_{17}| \leq (1+u)^4 \cdot 2^{-\min(\alpha_o + \beta_o + \beta_u, \alpha_o + \alpha_u + \beta_o)} \cdot |r_h|.$$

$$|t_{18}| \leq (\frac{1+u}{(1-u)^2} \cdot (u^2 - u^3 + (-u^3 + u) \cdot 2^{-\alpha_o - \beta_o})$$

$$+ (1+u)^4 \cdot 2^{-\min(\alpha_o + \beta_o + \beta_u, \alpha_o + \alpha_u + \beta_o)}) \cdot (1+u) \cdot |r_h|$$
As $\alpha_o + \beta_o < \alpha_o + \beta_o + \beta_u$ and $\alpha_o + \beta_o < \alpha_o + \beta_o + \alpha_u \Rightarrow$

$$|t_{18}| \leq (\frac{(1+u)^2}{(1-u)^2} \cdot (u^2 - u^3 + ((-u^3 + u) + (1+u)^3 \cdot (1-u^2)) \cdot 2^{-\alpha_o - \beta_o})) \cdot |r_h|$$

• $t_{19} + t_{20} = Add112(t_{14}, t_{18})$ according to $Add112 \Rightarrow t_{19} + t_{20} = t_{14} + t_{18}$ and $|t_{20}| \leq u.|t_{19}| \Rightarrow$, we simplify:

$$|t_{19}| \le (u + 2^{-\alpha_o - \beta_o + 1}).|r_h|$$

 $|t_{20}| \le u.(u + 2^{-\alpha_o - \beta_o + 1}).|r_h|$

• $t_{21} + t_{22} = Add222(t_2, t_3, t_4, t_5)$ according to $Add222 \Rightarrow t_{21} + t_{22} = (t_2 + t_3 + t_4 + t_5).(1 + \epsilon_2)$ with $|\epsilon_2| \leq 6.u^2$ and $|t_{22}| \leq 6.u.|t_{21}| \Rightarrow$

$$|t_{21}| \le \frac{(u+1)^2 \cdot (1+6u^2)}{(1-u)^2} \cdot (2^{-\alpha_o} + 2^{-\beta_o}) \cdot |r_h|$$

$$|t_{22}| \le \frac{u.(u+1)^2.(1+6u^2)}{(1-u)^2}.(2^{-\alpha_o}+2^{-\beta_o}).|r_h|$$

• $r_m + r_\ell = Add222(t_{21}, t_{22}, t_{19}, t_{20})$, according to $Add222 \Rightarrow |r_\ell| \leq 6.u.|r_m|$. After the calculations, we have:

$$|r_m| \le \frac{1 + 6.u^2}{1 - u^2} \cdot \left(\frac{(u+1)^3 \cdot (1 + 6u^2)}{(1 - u)^2} \cdot (2^{-\alpha_o} + 2^{-\beta_o}) + (u+1) \cdot (u + 2^{-\alpha_o - \beta_o + 1})\right) \cdot |r_h|$$

after simplification:

$$|r_m| \le .(u + 2^{-\alpha_o} + 2^{-\beta_o} + 2^{-\alpha_o - \beta_o + 1}).|r_h|$$

and

$$|r_{\ell}| \le u.(u + 2^{-\alpha_o} + 2^{-\beta_o} + 2^{-\alpha_o - \beta_o + 1}).|r_h|$$

Theorem 12 (Relative error algorithm Mul333) Let (a_h, a_m, a_ℓ) and (b_h, b_m, b_ℓ) are **triple-double** numbers and the arguments of the function Mul233.

So:

$$r_h + r_m + r_\ell = (a_h + a_m + a_\ell).(b_h + b_m + b_\ell).(1 + \epsilon)$$

with

$$|\epsilon| <$$

$$32.u^3 + 3.u \cdot (2^{-\beta_o + 2} + 2^{-\alpha_o + 2}) + 2^{-\beta_o - \beta_u + 3} + 2^{-\alpha_o - \alpha_u + 3} + 2^{-\alpha_o - \alpha_u - \beta_u + 3} + 2^{-\alpha_o - \beta_o - \beta_u + 3}$$

Proof According to Mul333, we have:

•
$$r_m + r_\ell = Add222(t_{21}, t_{22}, t_{19}, t_{20})$$
 based on $Add222 \Rightarrow r_m + r_\ell = (t_{21} + t_{22} + t_{19} + t_{20})(1 + \epsilon_1)$ with $|\epsilon_1| \le 6.u^2$

•
$$t_{21} + t_{22} = (t_2 + t_3 + t_4 + t_5).(1 + \epsilon_2)$$
 with $|\epsilon_2| \le 6.u^2$

•
$$t_2 + t_3 = a_h.b_m$$
 and $t_4 + t_5 = a_m.b_m \Rightarrow$

$$t_{21} + t_{22} = (a_h \cdot b_m + a_m \cdot b_m) \cdot (1 + \epsilon_2)$$

$$r_m + r_\ell = ((a_h \cdot b_m + a_m \cdot b_m) \cdot (1 + \epsilon_2) + t_{19} + t_{20})(1 + \epsilon_1)$$

- $t_{19} + t_{20} = t_{14} + t_{18}$ thanks to Add112
- $t_{18} = (t_{16} + t_{17})(1 + \epsilon_3)$ with $|\epsilon_3| \le u$ thanks to the collary 5
- $t_{17} = (t_{12} + t_{13})(1 + \epsilon_4)$ with $|\epsilon_4| \le u$ thanks to the collary 5
- $t_{12} = (t_8 + t_9)(1 + \epsilon_7)$ with $|\epsilon_7| \le u$ thanks to the collary 5
- $t_8 = (a_h.b_\ell).(1 + \epsilon_{11})$ with $|\epsilon_{11}| \leq u$ and $t_9 = (a_\ell.b_h).(1 + \epsilon_{10})$ with $|\epsilon_{10}| \leq u$ thanks to the collary 5

•

$$t_{12} = ((a_h.b_\ell).(1 + \epsilon_{11}) + (a_\ell.b_h).(1 + \epsilon_{10}))(1 + \epsilon_7)$$

$$t_{17} = (((a_h.b_\ell).(1 + \epsilon_{11}) + (a_\ell.b_h).(1 + \epsilon_{10}))(1 + \epsilon_7) + t_{13})(1 + \epsilon_4)$$

- $t_{13} = (t_{10} + t_{11}).(1 + \epsilon_6)$ with $|\epsilon_6| \le u$
- $t_{10} = a_m.b_{\ell}.(1+\epsilon_9)$ and $t_{11} = a_{\ell}.b_m.(1+\epsilon_8)$ with $|\epsilon_8| \le u$ and $|\epsilon_9| \le u$

$$t_{13} = (a_m.b_\ell.(1+\epsilon_9) + a_\ell.b_m.(1+\epsilon_8)).(1+\epsilon_6)$$

$$t_{17} = (((a_h.b_\ell).(1+\epsilon_{11}) + (a_\ell.b_h).(1+\epsilon_{10}))(1+\epsilon_7) + (a_m.b_\ell.(1+\epsilon_9) + a_\ell.b_m.(1+\epsilon_8)).(1+\epsilon_6))(1+\epsilon_7) + (a_\ell.b_\ell).(1+\epsilon_9) + (a_\ell.$$

At the end of the calculation of t_{17} , we have:

$$t_{17} = a_h . b_\ell . (1 + \epsilon_4 + \epsilon_7 + \epsilon_{11} + \epsilon_4 . \epsilon_7 + \epsilon_4 . \epsilon_{11} + \epsilon_7 . \epsilon_{11} + \epsilon_4 . \epsilon_7 . \epsilon_{11})$$

$$+ a_\ell . b_h . (1 + \epsilon_4 + \epsilon_7 + \epsilon_{10} + \epsilon_4 . \epsilon_7 + \epsilon_4 . \epsilon_{10} + \epsilon_7 . \epsilon_{10} + \epsilon_4 . \epsilon_7 . \epsilon_{10})$$

$$+ a_m . b_\ell . (1 + \epsilon_4 + \epsilon_6 + \epsilon_9 + \epsilon_4 . \epsilon_6 + \epsilon_4 . \epsilon_9 + \epsilon_6 . \epsilon_9 + \epsilon_4 . \epsilon_6 . \epsilon_8)$$

$$+ a_m . b_\ell . (1 + \epsilon_4 + \epsilon_6 + \epsilon_8 + \epsilon_4 . \epsilon_6 + \epsilon_4 . \epsilon_8 + \epsilon_6 . \epsilon_8 + \epsilon_4 . \epsilon_6 . \epsilon_8)$$

• $t_{16} = (t_7 + t_{15}).(1 + \epsilon_5)$ with $|\epsilon_5| \le u$

$$t_{18} = ((t_7 + t_{15}).(1 + \epsilon_5) + t_{17})(1 + \epsilon_3)$$

$$t_{19} + t_{20} = t_{14} + ((t_7 + t_{15}).(1 + \epsilon_5) + t_{17})(1 + \epsilon_3)$$

$$t_{19} + t_{20} = t_{14} + (t_7 + t_{15}).(1 + \epsilon_3 + \epsilon_5 + \epsilon_3.\epsilon_5) + t_{17}.(1 + \epsilon_3)$$

$$t_{19} + t_{20} = t_{14} + t_{15} + t_7 + (t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_3.\epsilon_5) + t_{17}.(1 + \epsilon_3)$$
 As $t_{14} + t_{15} = t_1 + t_6$
$$t_{19} + t_{20} = t_1 + t_6 + t_7 + (t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_3.\epsilon_5) + t_{17}.(1 + \epsilon_3)$$
 As $t_6 + t_7 = a_m.b_m$
$$t_{19} + t_{20} = t_1 + a_m.b_m + (t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_3.\epsilon_5) + t_{17}.(1 + \epsilon_3)$$

$$r_m + r_\ell = ((a_h.b_m + a_m.b_m).(1 + \epsilon_2) + t_1 + a_m.b_m + (t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_3.\epsilon_5) + t_{17}.(1 + \epsilon_3))(1 + \epsilon_1)$$

$$r_m + r_\ell = ((a_h.b_m + a_m.b_m).(1 + \epsilon_2 + \epsilon_3 + \epsilon_2.\epsilon_3) + t_1.(1 + \epsilon_1) + t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3)$$

$$+ a_m.b_m.(1 + \epsilon_1) + (t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_1.\epsilon_5 + \epsilon_3.\epsilon_5 + \epsilon_1.\epsilon_3.\epsilon_5)$$
 We add $r_h \Rightarrow :$
$$r_h + r_m + r_\ell = r_h + t_1 + (a_h.b_m + a_m.b_m).(1 + \epsilon_2 + \epsilon_3 + \epsilon_2.\epsilon_3) + t_1.\epsilon_1 + t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3)$$

$$+ a_m.b_m.(1 + \epsilon_1) + (t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_1.\epsilon_5 + \epsilon_3.\epsilon_5 + \epsilon_1.\epsilon_3.\epsilon_5)$$

$$r_h + r_m + r_\ell = (a_h + a_m).(b_h, b_m, b_\ell) + a_\ell.(b_h, b_m) + \delta$$
 with $\delta = t_1.\epsilon_1 + (a_h.b_m + a_m.b_h).(\epsilon_2 + \epsilon_3 + \epsilon_2.\epsilon_3) + a_m.b_m.\epsilon_1 + (t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_1.\epsilon_3 + \epsilon_1.\epsilon_5 + \epsilon_3.\epsilon_5 + \epsilon_1.\epsilon_3.\epsilon_5) + t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_m.b_\ell + a_\ell.b_m)$. We search the upper bound of $|\delta|$ in function of $|a_h.b_h|$.

$$|\delta| \leq |t_1.\epsilon_1| + |(a_h.b_m + a_m.b_h).(\epsilon_2 + \epsilon_3 + \epsilon_2.\epsilon_3)| + |a_m.b_m.\epsilon_1| + |(t_7 + t_{15}).(\epsilon_3 + \epsilon_5 + \epsilon_1.\epsilon_3 + \epsilon_1.\epsilon_5 + \epsilon_3.\epsilon_5 + \epsilon_1.\epsilon_3.\epsilon_5)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_m.b_\ell + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_m.b_\ell + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_m.b_\ell + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_m.b_\ell + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h)| + |t_{17}.(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.b_h)| + |t_{17}.(1 + \epsilon_1.b_h)|$$

We begin to calculate t_{17} . $(1 + \epsilon_1 + \epsilon_3 + \epsilon_1.\epsilon_3) - (a_h.b_\ell + a_\ell.b_h + a_m.b_\ell + a_\ell.b_m)$:

$$t_{17}.(1+\epsilon_1+\epsilon_3+\epsilon_1.\epsilon_3)-(a_h.b_\ell+a_\ell.b_h+a_m.b_\ell+a_\ell.b_m)=$$

$$a_h.b_\ell.((1+\epsilon_4+\epsilon_7+\epsilon_{11}+\epsilon_4.\epsilon_7+\epsilon_4.\epsilon_{11}+\epsilon_7.\epsilon_{11}+\epsilon_4.\epsilon_7.\epsilon_{11}).(1+\epsilon_1+\epsilon_3+\epsilon_1.\epsilon_3)-1)$$

$$+a_\ell.b_h.((1+\epsilon_4+\epsilon_7+\epsilon_{10}+\epsilon_4.\epsilon_7+\epsilon_4.\epsilon_{10}+\epsilon_7.\epsilon_{10}+\epsilon_4.\epsilon_7.\epsilon_{10}).(1+\epsilon_1+\epsilon_3+\epsilon_1.\epsilon_3)-1)$$

$$+a_m.b_\ell.((1+\epsilon_4+\epsilon_6+\epsilon_9+\epsilon_4.\epsilon_6+\epsilon_4.\epsilon_9+\epsilon_6.\epsilon_9+\epsilon_4.\epsilon_6.\epsilon_9).(1+\epsilon_1+\epsilon_3+\epsilon_1.\epsilon_3)-1)$$

$$+a_m.b_\ell.((1+\epsilon_4+\epsilon_6+\epsilon_8+\epsilon_4.\epsilon_6+\epsilon_4.\epsilon_8+\epsilon_6.\epsilon_8+\epsilon_4.\epsilon_6.\epsilon_8).(1+\epsilon_1+\epsilon_3+\epsilon_1.\epsilon_3)-1)$$

$$+a_m.b_\ell.((1+\epsilon_4+\epsilon_6+\epsilon_8+\epsilon_4.\epsilon_6+\epsilon_4.\epsilon_8+\epsilon_6.\epsilon_8+\epsilon_4.\epsilon_6.\epsilon_8).(1+\epsilon_1+\epsilon_3+\epsilon_1.\epsilon_3)-1)$$
As $|\epsilon_i| \leq 6.u^2$ with $i \in [1,2]$ and $|\epsilon_j| \leq u$ with $3 \leq j \leq 11$

We have so $|t_{17}.(1+\epsilon_1+\epsilon_3+\epsilon_1.\epsilon_3)-(a_h.b_\ell+a_\ell.b_h+a_m.b_\ell+a_\ell.b_m)|$
it's equal to

$$|a_h.b_\ell+a_\ell.b_h+a_m.b_\ell+a_\ell.b_m|.(6.u^6+24.u^5+37.u^4+23.u^3+9.u^2+6.u+1)$$

 \Rightarrow

$$|\delta| \le 6.u^2.|t_1| + (6.u^3 + 6.u^2 + u).|a_h.b_m + a_m.b_h| + 6.u^2.|a_m.b_m|$$
$$+(6.u^4 + 12.u^3 + u^2 + 2.u).|t_7 + t_{15}|$$
$$+(6.u^6 + 24.u^5 + 37.u^4 + 23.u^3 + 9.u^2 + 6.u + 1).|a_h.b_\ell + a_\ell.b_h + a_m.b_\ell + a_\ell.b_m|$$

. According to the conditions of Mul333:

$$|a_{h}.b_{m}| \leq 2^{-\beta_{o}}.|a_{h}.b_{h}|$$

$$|a_{h}.b_{\ell}| \leq 2^{-\beta_{o}-\beta_{u}}.|a_{h}.b_{h}|$$

$$|a_{m}.b_{h}| \leq 2^{-\alpha_{o}}.|a_{h}.b_{h}|$$

$$|a_{m}.b_{m}| \leq 2^{-\alpha_{o}-\beta_{o}}.|a_{h}.b_{h}|$$

$$|a_{m}.b_{\ell}| \leq 2^{-\alpha_{o}-\beta_{o}-\beta_{u}}.|a_{h}.b_{h}|$$

$$|a_{\ell}.b_{h}| \leq 2^{-\alpha_{o}-\alpha_{u}}.|a_{h}.b_{h}|$$

$$|a_{\ell}.b_{m}| \leq 2^{-\alpha_{o}-\alpha_{u}-\beta_{o}}.|a_{h}.b_{h}|$$

 \Rightarrow

$$|\delta| \leq 6.u^{2}.|t_{1}| + (6.u^{3} + 6.u^{2} + u).(2^{-\beta_{o}}.|a_{h}.b_{h}| + 2^{-\alpha_{o}}.|a_{h}.b_{h}|) + 6.u^{2}.2^{-\alpha_{o} - \beta_{o}}.|a_{h}.b_{h}|$$

$$+ (6.u^{4} + 12.u^{3} + u^{2} + 2.u).|t_{7} + t_{15}|$$

$$+ (6.u^{6} + 24.u^{5} + 37.u^{4} + 23.u^{3} + 9.u^{2} + 6.u + 1)$$

$$\times (2^{-\beta_{o} - \beta_{u}}.|a_{h}.b_{h}| + 2^{-\alpha_{o} - \alpha_{u}}.|a_{h}.b_{h}| + 2^{-\alpha_{o} - \beta_{o} - \beta_{u}}.|a_{h}.b_{h}| + 2^{-\alpha_{o} - \alpha_{u} - \beta_{o}}.|a_{h}.b_{h}|)$$

We seek the upper bound of $|t_1|$ in function of $|a_h.b_h|$. As $r_h + t_1 = a_h \cdot b_h$ thanks to $Add112 \Rightarrow |t_1| \leq u \cdot |r_h| \Rightarrow$

$$|r_h| \le \frac{1}{1-u}.|a_h.b_h|$$

 \Rightarrow

$$|t_1| \le \frac{u}{1-u}.|a_h.b_h|$$

Now, doing the same operation:

$$|t_7| \le \frac{u}{1-u}.|a_m.b_m|$$

$$|t_7| \le \frac{u}{1-u} \cdot 2^{-\alpha_o - \beta_o} \cdot |a_h \cdot b_h|$$

$$|t_{15}| \le \frac{u}{1-u}.|t_1+t_6|$$

 $|t_{15}| \le \frac{u}{1-u}.(|t_1|+|t_6|)$

We have:

$$|t_6| \le \frac{1}{1-u} \cdot 2^{-\alpha_o - \beta_o} \cdot |a_h \cdot b_h|$$

 \Rightarrow

$$|t_{15}| \le \frac{u}{1-u} \cdot \left(\frac{u}{1-u} \cdot |a_h \cdot b_h| + \frac{1}{1-u} \cdot 2^{-\alpha_o - \beta_o} \cdot |a_h \cdot b_h|\right)$$
$$|t_{15}| \le \frac{u}{(1-u)^2} \cdot (u + 2^{-\alpha_o - \beta_o}) \cdot |a_h \cdot b_h|$$

. After simplify, we have:

$$|\delta| \le (8.u^3 + 3.u \cdot (2^{-\beta_o} + 2^{-\alpha_o}) + 2^{-\beta_o - \beta_u + 1} + 2^{-\alpha_o - \alpha_u + 1} + 2^{-\alpha_o - \alpha_u - \beta_u + 1} + 2^{-\alpha_o - \beta_o - \beta_u + 1}) \cdot |a_h \cdot b_h|$$

We seek the upper bound of $|a_h.b_h|$ in function of $|a_m.(b_h,b_m,b_\ell) + a_\ell.(b_h,b_m) + a_h.(b_m,b_\ell)|$

$$|a_{m}.(b_{h},b_{m},b_{\ell}) + a_{\ell}.(b_{h},b_{m}) + a_{h}.(b_{m},b_{\ell})| \leq$$

$$(2^{-\beta_{o}} + 2^{-\beta_{o}-\beta_{u}} + 2^{-\alpha_{o}} + 2^{-\alpha_{o}-\beta_{o}-\beta_{u}+1} + 2^{-\alpha_{o}-\alpha_{u}+1} + 2^{-\alpha_{o}-\beta_{o}+1} + 2^{-\alpha_{o}-\alpha_{u}-\beta_{o}+1}).|a_{h}.b_{h}|$$
But $1 - (2^{-\beta_{o}} + 2^{-\beta_{o}-\beta_{u}} + 2^{-\alpha_{o}} + 2^{-\alpha_{o}-\beta_{o}-\beta_{u}} + 2^{-\alpha_{o}-\alpha_{u}} + 2^{-\alpha_{o}-\beta_{o}} + 2^{-\alpha_{o}-\beta_{o}})$

$$2^{-\alpha_{o}-\alpha_{u}-\beta_{o}} \geq \frac{1}{4} \text{ We have } :$$

$$|a_h.b_h| \le 4.|a_m.(b_h,b_m,b_\ell) + a_\ell.(b_h,b_m) + a_h.(b_m,b_\ell)|$$

So we have:

$$|\delta| \le$$

$$(32.u^3 + 3.u.(2^{-\beta_o + 2} + 2^{-\alpha_o + 2}) + 2^{-\beta_o - \beta_u + 3} + 2^{-\alpha_o - \alpha_u + 3} + 2^{-\alpha_o - \alpha_u - \beta_u + 3} + 2^{-\alpha_o - \beta_o - \beta_u + 3})$$

$$\times |(a_h + a_m + a_\ell).(b_h + b_m + b_\ell)|$$

so:

$$|\epsilon| \ge 32.u^3 + 3.u \cdot (2^{-\beta_o + 2} + 2^{-\alpha_o + 2}) + 2^{-\beta_o - \beta_u + 3} + 2^{-\alpha_o - \alpha_u + 3} + 2^{-\alpha_o - \alpha_u - \beta_u + 3} + 2^{-\alpha_o - \beta_o - \beta_u + 3}$$

Chapter 5

log

5.1 crlogfast

See algorithm 12

In practice with the C program, we found as an error bound $err_{fast} = 2^{58.4}$. (see the C program in appendix).

5.2 crlogaccurate

See algorithm 13

In practice with the C program, we found as an error bound $err_{accurate} = 2^{98.4}$. (see appendix the C program).

5.3 crlogadvanced

See algorithm 14

$5.4 \log$

See algorithm 15

The log function is composed of the 3 $cr \log$. We start using $cr \log_{fast}$, if this function cannot find the result of the log, it returns to $cr \log_{accurate}$, it does the same thing as the previous function. The last $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ of the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate all the calculations not solved by the previous $cr \log_{advanced}$ will calculate $cr \log_{advanced}$ will calcu

The relative errors calculated for $crlog_{fast}$ and $crlog_{accurate}$, will be used for the function to know if the calculation passes for each.

Let x be a **double** $,err_{fast}$ is the relative error calculated for $cr \log_{fast}$ and err_{acc} is that of $cr \log_{accurate}$.

The computation of $\log(x)$ proceeds by first computing $cr \log_{fast}(x)$ which gives us a **double-double**. We will name this **double-double** (h_1, ℓ_1) .

We take $right = h_1 + err_{fast} * h_1 + \ell_1$ and $left = h_1 - err_{fast} * h_1 + \ell_1$. If right = left then we have the result of $\log(x)$ otherwise we calculate with $cr \log_{accurate}(x)$. Let (h_2, ℓ_2) be the result of $cr \log_{accurate}(x)$.

We set $right = h_2 + err_{accurate} * h_2 + \ell_2$ and $left = h_2 - err_{accurate} * h_2 + \ell_2$. If right = left then we have the result of $\log(x)$ otherwise we calculate with $cr \log_{advancedaccurate}(x)$. At the end of the log function calculation, we have the real double value of $\log(x)$.

Conclusion

. In our report, we studied addition and multiplication functions which have as arguments either **Double** , **Double-Double** or **Triple-Double** that resulted in **Double-Double** or in **Triple-Double**.

We noticed that the calculations of **Double-Double** numbers were more precise than the calculations of **Doubles** numbers and also that the operations of **Triple-Double** numbers were more precise than those of **Double-Double** numbers and **Double** numbers.

We also shown how to implement the $cr \log$ and calculate the bounds of each thanks to their relative error calculations. We managed to implement our logarithm thanks to the bounds of $cr \log_{fast}$ and $cr \log_{accurate}$. Inaddition to this, We have shown how bounds were used in the logarithm.

For each $cr \log$, we calculated α_i and also their $\log(\alpha_i)$. $cr \log fast$ used α_i which are **Double** and have precision of 71 bits, while $cr \log accurate$ used α_i which are **Double-Double** and have a precision of approximately 107 bits and the $cr \log_{advanced}$ used α_i which are **Triple-Double** and which have a precision of about 160 bits.

We noticed that, while the $cr \log_{fast}$ has faster internal calculations than those of $cr \log_{accurate}$, the internal operations of $cr \log_{accurate}$ is faster than those of du $cr \log_{advanced}$.

In our research, we managed to implement the logarithm with the calculation algorithms, and We tested it on millions of worst cases.

To our mind , it would be necessary to test the speed of the code, by trying to make it faster. Like for example, changing the $cr \log_{fast}$ so that its error bound would be more smaller.

Thanks to this logarithm, we could implement the logarithm in base 2 (\log_2) and base 10 (\log_{10}). Moreover, we could Implement the logarithm with correct rounding for binary80 and binary128 and the same for (\log_2)

and (\log_{10}) .

Annexes

In appendix, We find each algorithm with their programs in **Sage** and in **C**.

Add112

Algorithm 1 Algorithm Add112 (FastTwoSum)

Input: a and b are 53-bit floating-point numbers

Condition: $|a| \ge |b|$

Output: s and t are 53-bit floating-point numbers : s: main value and t: error value.

```
1: s = a + b
2: z = s - a
3: t = b - z
4: return s, t
```

```
1 #Input: a and b are 53-bit floating-point numbers
2 def Add112(a,b):
      \mathtt{s} = \mathtt{a} + \mathtt{b}
      z = s-a
4
      \mathtt{t} = \mathtt{b} {-} \mathtt{z}
       \# output: s and t are 53-bit floating-point numbers.
       \# s: main value and t: error value
       return s,t
1 /*a and b are double numbers */
2 void Add112(double a, double b, double *s, double *t){
      *{\tt s}={\tt a}+{\tt b}\;;
       double z = *s-a;
       *t = b-z;
5
       /* (s,t) is a double—double number*/
```

Add122

Algorithm 2 Algorithm Add122

```
Input: a is 53-bit floating-point numbers, b_h: main value and b_\ell: error value

Condition: |a| \geq |b_h|

Condition: |b_\ell| \leq u.|b_h|

Output: s and t are 53-bit floating-point numbers : s: main value and t: error value.

1: s, \ell = Add112(a, b_h)

2: t = \ell + b_\ell

3: return s, t
```

```
1 #Input : a is 53-bit floating-point numbers
2 \# bh: main value and bl: error value
3 def Add122(a,bh,bl):
      s,l = Add112(a,bh)
      \mathbf{t}=\mathbf{1}\mathbf{+bl}
      \# output: s and t are 53-bit floating-point numbers.
7
      \# s: main value and t: error value
      return s,t
1 /* a is double—number and (bh,bl) is double—double number*/
2 void Add122(double a, double bh, double bl, double *s, double *t){
      double 1;
      Add112(a, bh,s,&1);
4
      *{\tt t} = 1 + {\tt bl};
      /*(s,t) is a double-double number*/
```

Add222

Algorithm 3 Algorithm Add222 Input: a_h and b_h are main values, a_ℓ and b_ℓ are error values Condition: $|a_h| \ge |b_h|$ Condition: $|a_\ell| \le u.|a_h|$ and $|b_\ell| \le u.|b_h|$ Output: s and t are 53-bit floating-point numbers : s: main value and t: error value. 1: $s, \ell = Add112(a_h, b_h)$ 2: $m = \ell + a_\ell$ 3: $t = m + b_\ell$ 4: return s, t

```
1 #Input : ah and bh are main values
2 \# al and bl are error values
3 def Add222(ah,al,bh,bl):
       s,l = Add112(ah,bh)
       {\tt m} = {\tt l} {+} {\tt al}
5
6
       \mathtt{t}=\mathtt{m}{+}\mathtt{bl}
       \# output: s and t are 53-\mathrm{bit} floating-point numbers.
       \# s: main value and t: error value
       return s,t
1 /* (ah,al) and (bh,bl) double—double numbers*/
2 void Add222(double ah, double al, double bh, double bl, double *s, double *t){
       double l,m;
4
       Add112(ah,bh,s,\&1);
5
       \mathtt{m}=\mathtt{l}+\mathtt{al};
       *{\tt t}={\tt m}{+}{\tt bl};
       /*(s,t) is a double—double number*/
```

Mul112

Algorithm 4 Algorithm Mul112 (DEKKER-PRODUCT)

Input: a and b are 53-bit floating-point numbers.

Output: r_1 and r_2 are 53-bit floating-point numbers: r_1 : main value and r_2 : error value.

```
1: r1 = a \times b
2: r2 = FMA(a, b, -r1)
3: return (r_1, r_2)
```

```
1 /*a and b are double numbers */
2 void Mul112(double a, double b, double *r1, double *r2){
3     *r1 = a * b;
4     *r2 = __builtin_fma (a, b, -*r1);
5     /* (r1,r2) is a double number */
6 }

1 #Input : a and b are 53-bit floating-point numbers
2 def Mul112(a,b):
3     r1 = a*b
4     r2 = fma(a,b,-r1)
5     # output: r1 and r2 are 53-bit floating-point numbers.
6     # r1: main value and r2: error value
7     return (r1,r2)
```

Mul122

Algorithm 5 Algorithm Mul122 Input: a is 53-bit floating-point numbers, b_h : main value and b_ℓ : error value Condition: $|b_\ell| \le u.|b_h|$ Output: r_1 and r_2 are 53-bit floating-point numbers : r_1 : main value and r_2 : error value

1: $t_1, t_2 = Mul112(a, b_h)$ 2: $t_3 = a \times b_\ell$ 3: $t_4 = t_2 + t_3$ 4: $r_1, r_2 = Add112(t_1, t_4)$ 5: return (r_1, r_2)

```
1~\# Input: a is 53-bit floating-point numbers
 2 \# bh: main value and bl: error value
 3 def Mul122(a,bh,bl):
        t1,t2 = Mull112(a,bh)
 4
 5
        {\tt t3} = {\tt a*bl}
        \mathtt{t4}=\mathtt{t2}{+}\mathtt{t3}
 6
        r1,r2 = Add112(t1,t4)
        # output: r1 and r2 are 53-bit floating-point numbers.
 8
        \# r1: main value and r2: error value
10
        \underline{return}\ (\mathtt{r1},\!\mathtt{r2})
 1 /* a is double—number and (bh,bl) is double—double number*/
 2 void Mul122(double a, double bh, double bl, double *r1, double *r2){
        double t1,t2;
        Mul112(a,bh,&t1,&t2);
 4
        double t3,t4;
 5
        {\tt t3} = {\tt a*bl};
 7
        \mathtt{t4}=\mathtt{t2}\mathtt{+t3};
        Add112(t1,t4,r1,r2);
        /* (r1,r2) is a double number */
```

Mul222

Algorithm 6 Algorithm Mul222

Input: a_h and b_h are main values, a_ℓ and b_ℓ are error values

Output: r_1 and r_2 are 53-bit floating-point numbers: r_1 : main value and r_2 : error value.

```
1: t_1, t_2 = Mul112(a_h, b_h)

2: t_3 = a_h \times b_\ell

3: t_4 = b_h \times a_\ell

4: t_5 = t_3 + t_4

5: t_6 = t_2 + t_5

6: r_1, r_2 = Add112(t_1, t_6)

7: return (r_1, r_2)
```

```
1 #Input: ah and bh are main values
 2 \# al and bl are error values
 3 def Mul222(ah,al,bh,bl):
        t1,t2 = Mull112(ah,bh)
        \mathtt{t3} = \mathtt{ah} * \mathtt{bl}
        {\tt t4} = {\tt bh*al}
       \mathtt{t5}, \mathtt{t6} = \mathtt{Add112Cond}(\mathtt{t3}, \mathtt{t4})
        r1,r2 = Add222Cond(t1,t2,t5,t6)
        \# output: r1 and r2 are 53-\text{bit floating}-\text{point numbers}.
10
        \# r1: main value and r2: error value
        return (r1,r2)
11
 1 /* (ah,al) and (bh,bl) double—double numbers*/
 2 void Mul222(double ah, double al, double bh, double bl, double *r1, double *r2){
        double t1,t2;
        Mul112(ah, bh, &t1, &t2);
        double t3,t4;
        t3 = ah*b1;
        t4 = bh*al;
 7
        double t5,t6;
 9
        t5 = t3+t4;
        t6 = t2+t5;
        Add112(t1,t6,r1,r2);
11
        /* (r1,r2) is a double number */
13 }
```

Add133

Algorithm 7 Algorithm Add133

```
Input: a is a double number and b_h, b_m, b_\ell is a triple-double numbers. Condition: |a| \geq |b_h|, |b_m| \leq u.|b_h|, |b_\ell| \leq u.|b_m| and |b_\ell| \leq u^2.|b_h| Output: r_h, r_m, r_\ell is a triple-double numbers.

1: r_h, t_1 = Add112(a, b_h)
2: t_2, t_3 = Add112(t_1, b_m)
3: t_4 = t_3 + b_\ell
4: r_m, r_\ell = Add112(t_2, t_4)
5: return r_h, r_m, r_\ell
```

```
1 \# Input : a is a double numbers
 2 # bh,bm,bl is a triple—double numbers
 3 def Add133(a,bh,bm,bl):
       \mathtt{rh},\mathtt{t1} = \mathtt{Add112}(\mathtt{a},\mathtt{bh})
       t2,t3 = Add112(t1,bm)
       \mathtt{t4} = \mathtt{t3} + \mathtt{b1}
       rm,rl = Add112(t2,t4)
        #output :rh,rm,rl is a triple-double numbers
       return rh,rm,rl
 1 /* a is a double numbers
      bh,bm,bl is a triple—double numbers*/
 3 void Add133(double a, double bh, double bm, double bl, double *rh, double *rm, double *rl){
       double t1,t2,t3,t4;
       Add112(a,bh,rh,&t1);
       Add112(t1,bm,&t2,&t3);
 6
 7
       t4 = t3+b1;
       Add112(t2,t4,rm,rl);
        /*rh,rm,rl is a triple—double numbers*/
10 }
```

5

6

3

5

7

8

9

12 }

Add333

Algorithm 8 Algorithm Add333

Input: a_h , a_m , a_ℓ and b_h , b_m , b_ℓ are triple-double numbers Condition: $|b_h| \leq \frac{3}{4}.|a_h|$ Condition: $|a_m| \leq 2^{-\alpha_o} . |a_h|$ Condition: $|a_{\ell}| \leq 2^{-\alpha_u} . |a_m|$ Condition: $|b_m| \leq 2^{-\beta_o} . |b_h|$ Condition: $|b_{\ell}| \leq 2^{-\beta_u}.|b_m|$ Condition: $\alpha_o \ge 4$, $\alpha_u \ge 1$, $\beta_o \ge 4$, $\beta_u \ge 1$ **Output:** r_h , r_m , r_ℓ is a triple-double numbers. 1: $r_h, t_1 = Add112(a_h, b_h)$ 2: $t_2, t_3 = Add112Cond(a_m, b_m)$ 3: $t_7, t_4 = Add112(t_1, t_2)$ 4: $t_6 = a_\ell + b_\ell$ 5: $t_5 = t_3 + t_4$ 6: $t_8 = t_5 + t_6$ 7: $r_m, r_\ell = Add112(t_7, t_8)$ 8: return r_h , r_m , r_ℓ 1 def Add333(ah,am,al,bh,bm,bl): rh,t1 = Add112(ah,bh)t2,t3 = Add112(am,bm)t7,t4 = Add112(t1,t2) ${ t t6} = { t al} + { t bl}$ t5 = t3+t4t8 = t5 + t6rm,rl = Add112(t7,t8)#output :rh,rm,rl is a triple-double numbers 9 10 return rh,rm,rl 1 /*ah,am,al and bh,bm,bl are triple—double numbers*/ 2 void Add333(double ah, double am, double al, double bh, double bm, double bl, double *rh, double *rm, double *rl){ double t1,t2,t3,t4,t5,t6,t7,t8; Add112(ah,bh,rh,&t1); Add112Cond(am,bm,&t2,&t3); Add112(t1,t2,&t7,&t4);t6 = al+bl;t5 = t3+t4;t8 = t5 + t6;Add112(t7,t8,rm,rl); 10 /*rh,rm,rl is a triple—double numbers*/ 11

Mul133

Algorithm 9 Algorithm Mul133 Input: a is a double number and b_h , b_m , b_ℓ is a triple-double numbers. Condition: $|b_m| \leq 2^{-\beta_o}.|b_h|$ with $\beta_o \geq 2$ Condition: $|b_\ell| \leq 2^{-\beta_u}.|b_m|$ with $\beta_u \geq 2$ Output: r_h , r_m , r_ℓ is a triple-double numbers. 1: r_h , $t_2 = Mul112(a, b_h)$ 2: t_3 , $t_4 = Mul112(a, b_m)$ 3: $t_5 = a \times b_\ell$ 4: t_9 , $t_7 = Add112(t_2, t_3)$ 5: $t_8 = t_4 + t_5$ 6: $t_{10} = t_7 + t_8$ 7: r_m , $r_\ell = Add112(t_9, t_{10})$ 8: return r_h , r_m , r_ℓ

```
1 #Input : a is a double numbers
 2 # bh,bm,bl is a triple—double numbers
 3 def Mul133(a,bh,bm,bl):
       rh,t2 = Mull112(a,bh)
 4
       {\tt t3,t4} = {\tt Mull12(a,bm)}
 6
       {\tt t5} = {\tt a*bl}
       t9,t7 = Add112Cond(t2,t3)
       \mathtt{t8}=\mathtt{t4}{+}\mathtt{t5}
 8
9
       \mathtt{t10} = \mathtt{t7} {+} \mathtt{t8}
       rm,rl = Add112Cond(t9,t10)
10
        #output :rh,rm,rl is a triple-double numbers
11
12
       return rh,rm,rl
 1 /* a is a double numbers
      bh,bm,bl is a triple—double numbers*/
 3 void Mul133(double a, double bh, double bm, double bl, double *rh, double *rm, double *rl){
       double t2,t3,t4,t5,t7,t8,t9,t10;
       Mul112(a,bh,rh,&t2);
 5
       Mul112(a,bm,&t3,&t4);
       t5 = a*b1:
       Add112(t2,t3,&t9,&t7);
9
       t8 = t4+t5;
       t10 = t7 + t8;
10
11
       Add112(t9,t10,rm,rl);
12
        /*rh,rm,rl is a triple—double numbers*/
13 }
```

Mul233

```
Algorithm 10 Algorithm Mul233
   Input: a_h, a_\ell is a double-double and b_h, b_m, b_\ell is a triple-double numbers.
   Condition: |a_{\ell}| \leq u.|a_h|,
   Condition: |b_m| \leq 2^{-\beta_o} |b_h| with \beta_o \geq 2
   Condition: |b_{\ell}| \leq 2^{-\beta_u} |b_m| with \beta_u \geq 1.
   Output: r_h, r_m, r_\ell is a triple-double numbers.
     1: r_h, t_1 = Mul112(a_h, b_h)
     2: t_2, t_3 = Mul112(a_h, b_m)
     3: t_4, t_5 = Mul112(a_h, b_\ell)
     4: t_6, t_7 = Mul112(a_\ell, b_h)
     5: t_8, t_9 = Mul112(a_\ell, b_m)
     6: t_{10} = a_{\ell} \times b_{\ell}
     7: t_{11}, t_{12} = Add222(t_2, t_3, t_4, t_5)
     8: t_{13}, t_{14} = Add222(t_6, t_7, t_8, t_9)
     9: t_{15}, t_{16} = Add222(t_{11}, t_{12}, t_{13}, t_{14})
    10: t_{17}, t_{18} = Add112(t_1, t_{10})
    11: r_m, r_\ell = Add222(t_{17}, t_{18}, t_{15}, t_{16})
    12: return r_h, r_m, r_\ell
 1 #Input : ah,al is a double—double numbers
 2 # bh,bm,bl is a triple—double numbers
 3 def Mul233(ah,al,bh,bm,bl):
       rh,t1 = Mull112(ah,bh)
       \mathtt{t2},\mathtt{t3} = \mathtt{Mul112}(\mathtt{ah},\mathtt{bm})
 6
       t4,t5 = Mull112(ah,bl)
       {\tt t6,t7} = {\tt Mulll2(bh,al)}
 7
       t8,t9 = Mull112(al,bm)
 8
9
       t10 = al*bl
10
       t11,t12 = Add222Cond(t2,t3,t4,t5)
11
       t13,t14 = Add222Cond(t6,t7,t8,t9)
       t15,t16 = Add222Cond(t11,t12,t13,t14)
12
       t17,t18 = Add112Cond(t1,t10)
13
       rm,rl = Add222Cond(t17,t18,t15,t16)
14
        #output :rh,rm,rl is a triple-double numbers
15
16
       return rh,rm,rl
 1 /* ah, al is a double—double numbers
      bh,bm,bl is a triple—double numbers*/
 3 void Mul233(double ah, double al, double bh, double bm, double bl, double *rh, double *rm,
        double *rl){
```

double t1,t2,t3,t4,t5,t6,t7,t8,t9,t10; double t11,t12,t13,t14,t15,t16,t17,t18;

```
\begin{array}{l} {\tt Mul112(ah,bh,rh,\&t1);} \\ {\tt Mul112(ah,bm,\&t2,\&t3);} \end{array}
 7
            Mul112(ah,bl,&t4,&t5);
Mul112(bh,al,&t6,&t7);
 8
 9
            Mul112(al,bm,&t8,&t9);
t10 = al*bl;
10
11
            \begin{array}{l} \texttt{Add222(t2,t3,t4,t5,\&t11,\&t12);} \\ \texttt{Add222(t6,t7,t8,t9,\&t13,\&t14);} \end{array}
12
13
             Add222(t11,t12,t13,t14,&t15,&t16);
14
            Add112(t1,t10,&t17,&t18);
Add222(t17,t18,t15,t16,rm,rl);
15
16
             /*rh,rm,rl is a triple—double numbers*/
17
18 }
```

Mul333

Algorithm 11 Algorithm Mul333 **Input:** a_h , a_m , a_ℓ and b_h , b_m , b_ℓ are triple-double numbers Condition: $|a_m| \leq 2^{-\alpha_o} . |a_h|$ Condition: $|a_{\ell}| \leq 2^{-\alpha_u}.|a_m|$ Condition: $|b_m| \leq 2^{-\beta_o} |b_h|$ Condition: $|b_{\ell}| \leq 2^{-\beta_u}.|b_m|$ Condition: $\alpha_o \geq 2$, $\alpha_u \geq 2$, $\beta_o \geq 2$, $\beta_u \geq 2$ **Output:** r_h , r_m , r_ℓ is a triple-double numbers 1: $r_h, t_1 = Mul112(a_h, b_h)$ 2: $t_2, t_3 = Mul112(a_h, b_m)$ 3: $t_4, t_5 = Mul112(a_m, b_h)$ 4: $t_6, t_7 = Mul112(a_m, b_m)$ 5: $t_8 = a_h \times b_\ell$ 6: $t_9 = a_\ell \times b_h$ 7: $t_{10} = a_m \times b_\ell$ 8: $t_{11} = a_{\ell} \times b_m$ 9: $t_{12} = t_8 + t_9$ 10: $t_{13} = t_{10} + t_{11}$ 11: $t_{14}, t_{15} = Add112(t_1, t_6)$ 12: $t_{16} = t_7 + t_{15}$ 13: $t_{17} = t_{12} + t_{13}$ 14: $t_{18} = t_{16} + t_{17}$ 15: $t_{19}, t_{20} = Add112(t_{14}, t_{18})$ 16: $t_{21}, t_{22} = Add222(t_2, t_3, t_4, t_5)$ 17: $r_m, r_\ell = Add222(t_{21}, t_{22}, t_{19}, t_{20})$ 18: return r_h , r_m , r_ℓ

```
1 #Input : ah,am,al and bh,bm,bl are triple—double numbers
 2 def Mul333(ah,am,al,bh,bm,bl):
          rh,t1 = Mull112Cond(ah,bh)
          {\tt t2,t3} = {\tt Mull112Cond}({\tt ah,bm})
          t4,t5 = Mull112Cond(bh,am)
 6
          t6,t7 = Mull112Cond(am,bm)
          {\tt t8} = {\tt ah*bl}
          {	t t9} = {	t al}*{	t bh}
 8
          \mathtt{t10} = \mathtt{am*bl}
 9
10
          \mathtt{t11} = \mathtt{al*bm}
11
          t12 = t8 + t9
12
           \mathtt{t13} = \mathtt{t10} {+} \mathtt{t11}
          \mathtt{t14}, \mathtt{t15} = \mathtt{Add112Cond}(\mathtt{t1}, \mathtt{t6})
13
          \mathtt{t16} = \mathtt{t7} {+} \mathtt{t15}
```

```
15
       \mathtt{t17} = \mathtt{t12} {+} \mathtt{t13}
        \mathtt{t18} = \mathtt{t16} {+} \mathtt{t17}
16
17
        t19,t20 = Add112Cond(t14,t18)
18
        \mathtt{t21},\mathtt{t22} = \mathtt{Add222Cond}(\mathtt{t2},\mathtt{t3},\mathtt{t4},\mathtt{t5})
       rm,rl = Add222Cond(t21,t22,t19,t20)
19
20
        #output :rh,rm,rl is a triple—double numbers
21
        return rh,rm,rl
 1\ /*ah,am,al and bh,bm,bl are triple—double numbers
*/
 2 void Mul333(double ah, double am, double al, double bh, double bm, double bl, double *rh,
        double *rm, double *rl){
        double t1,t2,t3,t4,t5,t6,t7,t8,t9,t10;
 3
 4
        double t11,t12,t13,t14,t15,t16,t17,t18,t19,t20;
        double t21,t22;
 5
 6
        Mul112(ah,bh,rh,&t1);
        Mul112(ah,bm,&t2,&t3);
 7
 8
        Mul112(am,bh,&t4,&t5);
 9
       Mul112(am,bm,&t6,&t7);
        t8 = ah*b1;
10
11
        {\tt t9} = {\tt al*bh};
       t10 = am*b1;
12
13
        t11 = al*bm;
        t12 = t8+t9;
14
        t13 = t10+t11;
        Add112(t1,t6,&t14,&t15);
16
17
        t16 = t7 + t15;
        t17 = t12+t13;
18
        t18 = t16 + t17;
19
20
        Add112(t14,t18,&t19,&t20);
21
        Add222(t4,t5,t2,t3,&t21,&t22);
22
        Add222(t21,t22,t19,t20,rm,rl);
23
        /*rh,rm,rl is a triple—double numbers*/
24 }
```

crlogfast

```
Algorithm 12 Algorithm crlog_{fast} for x is not to close to 1
Condition: x is a double number and is not to close to 1.
Condition: m_1 and i are results made after calculating x.
Condition: table_{\alpha_i} \ table_{\log(\alpha_i)} and are calculated with our method with
     Tang and Gal's method.
Condition: log(2) in double-double: h_{log(2)} = 0x1.62e42fefa38p - 1 and
     l_{\log 2} = 0x1.ef35793c7673p - 45
Input: m_1 is a double number, i is a Integer number. and table_{\alpha_i} and
     table_{\log(\alpha_i)} are tables of double numbers and f a polynomial of degree
     with coefficient f_i calculated with Sollya with 1 \le j \le 7
Condition: 1 \le m_1 \le 2, \ 0 \le i < 256 and size of table_{\alpha_i} = 256 and size of
     table_{\log(\alpha_i)} = 256, \alpha_i = table_{\alpha_i}[i] and \log_{\alpha_i} = table_{\log(\alpha_i)}[i]
Output: \log(x) in double number.
 1: h_r, \ell_r = Mul112(m_1, \alpha_i)
 2: h, \ell = Add122(-1, h_r, \ell_r)
 3: hh_2 = h \times h
 4: ff_5 = f_5 \times h
 5: ff_4 = f_4 + ff_5
 6: ff_7 = f_7 \times h
 7: ff_6 = f_6 + ff_7
 8: cr_{logfast_2} = hh_2 \times ff_6
 9: cr_{logfast_1} = ff_4 + cr_{logfast_2}
10: hh_2, \ell\ell_2 = Mul222(h, \ell, h, \ell)
11: hh_4, \ell\ell_4 = Mul222(hh_2, \ell\ell_2, hh_2, \ell\ell_2)
12: ffh_0, ff\ell_0 = Mul122(f_1, h, \ell)
13: ffh_3, ff\ell_3 = Mul122(f_3, h, \ell)
14: ffh_2, ff\ell_2 = Add122(f_2, ffh_3, ff\ell_3)
15: fhx_2, ff\ell x_2 = Mul222(hh_2, \ell\ell_2, ffh_2, ff\ell_2)
16: ffhx_4, ff\ell x_4 = Mul122(cr_{logfast_1}, hh_4, \ell\ell_4)
17: ffhx_0, ff\ell x_0 = Add222(ffh_0, ff\ell_0, ffhx_2, ff\ell x_2)
18: h_4, \ell_4 = Add222(ffhx_0, ff\ell x_0, ffhx_4, ff\ell x_4)
19: h_5, \ell_5 = Add122(\log_{\alpha_i}, h_4, \ell_4)
20: e_{hlog2}, e_{\ell log2} = Mul122(e, h_{log2}, \ell_{log2})
21: h_6, \ell_6 = Add222(e_{hlog2}, e_{\ell log2}, h_5, \ell_5)
22: return (h_6, \ell_6)
```

```
1 # Input : x is a 53-bit floating-point numbers
 2 # f is an approximation polynomial of log(1+x)
 3 def cr_log_fast_path(x,f):
       \# \mathtt{If} \ \mathtt{x} \ \mathtt{is} \ \mathtt{NAN} \ \underline{\mathtt{return}} \ \mathtt{NAN}
 5
       if x == NaN:
 6
 7
            return NaN
 8
 9
       \# If x=+0 or -0 return -INFINITY
10
       if x == 0:
           return -oo
11
12
       \# If x = +INFINITY \frac{return}{return} +INFINITY
13
14
       if x == +\infty:
15
           return +oo
16
17
       \# if x<1 and close to 1 return with the function f
       if (x > RR(0x1.fe7814e49392fp-1.16) and x < 1):
18
19
            # we calculated len(f.list())
20
21
            F =f.list()
22
            F.reverse()
            lenf = len(F)
23
24
25
            hr, lr = Add112(-1,x)
26
27
            \# F[i] is the coefficient of each monomial :x^{(len(F)-1-i)}
28
            \# We are turning F[0] in hf,lf with hf,: main value
29
30
            \# lf: error value
            hf,lf = Split(F[0])
31
32
            \# multiply (hr,lr) by (hf,lf) result (h,l)
33
34
            h,l = Mul222(hf,lf,hr,lr)
35
36
            for i in range(1,lenf-1):
37
                #We are turning F[i] in hfi,lfi with hfi: main value
                \# lfi: error value
38
39
                hfi,lfi = Split(F[i])
40
                \# add (hfi,lfi) by (h,l) result (h1,l1)
41
42
                h1,l1 = Add222(hfi,lfi,h,l)
43
                \# multiply (hr,lr) by (h1,l1) result (h,l)
44
                h,1 = Mul222(h1,l1,hr,lr)
45
46
47
            return h,1
48
49
       \# if x>1 and close to 1 return with function f
50
       if x > 1 and x < RR(0x1.00068db8bac71p+0.16):
            \# we calculated len(f.list())
51
52
            F =f.list()
           F.reverse()
53
            lenf = len(F)
54
55
56
            hr, lr = Add112(x, -1)
57
            \# F[i] is the coefficient of each monomial :x^{(len(F)-1-i)}
58
59
```

```
60
              #We are turning F[0] in hf,lf with hf,: main value
              # lf: error value
 61
 62
              hf,lf = Split(F[0])
 63
              \# multiply (hr,lr) by (hf,lf) result (h,l)
 64
              h,l = Mul222(hf,lf,hr,lr)
 65
 66
 67
              for i in range(1,lenf-1):
                   #We are turning F[i] in hfi,lfi with hfi: main value
 68
 69
                   # lfi: error value
                   hfi,lfi = Split(F[i])
 70
 71
                   \# add (hfi,lfi) by (h,l) result (h1,l1)
 72
 73
                   h1,l1 = Add222(hfi,lfi,h,l)
 74
 75
                   \# multiply (hr,lr) by (h1,l1) result (h,l)
                   \mathbf{h}, \mathbf{l} = \mathtt{Mul222}(\mathbf{h1}, \mathbf{l1}, \mathbf{hr}, \mathbf{lr})
 76
 77
 78
              return h,1
 79
 80
 81
          \# s represents the sign
 82
          #e represents the exposant
 83
          \# m represents fraction
 84
          (s,m,e) = RR(x).sign_mantissa_exponent()
          \#e = e - 53
 85
 86
 87
 88
         \mathbf{e} = \mathbf{e}{+}53
 89
 90
 91
          #If x is a negative number return NAN
         if s==-1 and e!=0:
 92
 93
              return NaN
 94
 95
          \# If x is a subnormal
 96
         if s==1 and e<0 and m!=0:
 97
 98
              \mathtt{v}=\mathtt{m}
              \mathbf{e} = \mathbf{e}{-}1
 99
100
              while v < 2^52:
101
102
                   v*=2
103
                   \mathbf{e} = \mathbf{e} - 1
104
              m1 = v*2.^(-52)
105
106
107
          \# \mathtt{If}\ \mathtt{x}\ \mathtt{is}\ \mathtt{normal}
108
109
110
              \#print(s,e,m)
111
              m1 = m*2.^(-52)
              \mathtt{e}=\mathtt{e}{-1}
112
113
114
115
116
117
         binary = RR(m1).str(2)[2:10] \# we recover the 8 bits after the initial 1.
118
         i = int(binary, 2) \# i  is the 8-bit integer
```

```
119
120
121
        {\tt halpha\_i\_m,lalpha\_i\_m} = {\tt Split}({\tt RR}({\tt table\_alpha\_modified[i]},16))~\#{\tt table}~{\tt computed}~{\tt for}
              all i of alpha_i_m
122
                                            # such that log(alpha_i_m) is 71 bits accurate.
123
        hlog_alpha_i_m,llog_alpha_i_m= Split(RR(table_log_alpha_modified[i],16))
124
125
        # multiply (halpha_i_m,lalpha_i_m) by (hm,lm)
126
        hr,lr = Mul122(m1,halpha_i_m,lalpha_i_m)
127
        \# add hR,1R by (-1) in double,double
128
129
        h,l = Add122(-1,hr,lr)
130
131
        # we calculated len(f.list())
132
        F =f.list()
133
        F.reverse()
134
        lenf = len(F)
135
136
        \# F[i] is the coefficient of each monomial :x^{(len(F)-1-i)}
137
138
        #We are turning F[0] in hf,lf with hf,: main value
        \# lf: error value
139
        hf,lf = Split(F[0])
140
141
142
        \# multiply (hr,lr) by (hf,lf) result (h,l)
143
        h1,l1 = Mul222(hf,lf,h,l)
144
145
        for i in range(1,lenf-1):
            \# We are turning F[i] in hfi,lfi with hfi: main value
146
147
             # lfi: error value
            hfi,lfi = Split(F[i])
148
149
150
             \# add (hfi,lfi) by (h1,l1) result (h2,l2)
            h2,12 = Add222(hfi,lfi,h1,l1)
151
152
            \# multiply (h,1) by (h2,12) result (h1,11)
153
154
            h1,l1 = Mul222(h2,l2,h,l)
155
156
157
        # Add (h1,l1) by (hlog_alpha_i_m,llog_alpha_i_m)
158
159
        \verb|h5,15| = \verb|Add222(hlog_alpha_i_m,llog_alpha_i_m,h1,l1)|
160
161
162
        \# \log(2) = (h_{\log 2}, l_{\log 2})
163
        h_{\log 2} = RR(0x1.62e42fefa38p-1,16)
164
        1_{\log 2} = RR(0x1.ef35793c7673p-45,16)
165
166
        \#e*log(2)
        elog2_h,elog2_l = Mul122(RR(e),h_log2,l_log2)
167
168
169
        hlog_fi, llog_fi = Add222(elog2_h, elog2_l, h5, l5)
170
171
172
173
        \# \text{output: hlog\_fi and llog\_fi are } 53-\text{bit floating-point numbers.}
174
        # hlog_fi: main value and llog_fi: error value
175
176
        return RR(hlog_fi),RR(llog_fi)
```

```
1 static double cr_log_fast_path(double x, double *h6, double *l6){
       int s:
3
       int e;
 4
       uint64_t m;
       extract (x,&s,&e, &m);
5
 6
       /* If x is a negative number return NAN*/
 7
       if ((s == 1) \&\& (e != 0)){
 8
9
           return NAN;
10
        *If x is NAN return NAN*/
11
12
       if ((s == 0) \&\& (e == 0x7ff) \&\& (m!= 0))
           return NAN;
13
14
       /*If x=+0 ou -0 return INFINITY*/
15
       if (((s==1) \&\& (e==0) \&\& (m==0)) || ((s==0) \&\& (e==0) \&\& (m==0)))
16
           return - (0x1p1023 + 0x1p1023);
17
18
        /*If x = +INFINITY return +INFINITY */
19
       \inf ((s == 0) \&\& (e == 0x7ff)){
20
           return 0x1p1023 + 0x1p1023;
21
22
       /*The coefficients of the approximation polynomial of degree 7. */
23
       double f7 = 0x1.2152a2de69894p-3;
24
       double f6 = -0x1.555147415c204p-3;
25
26
       double f5 = 0x1.999997342c184p-3;
27
       double f3 = 0x1.5555555554cep-2;
28
29
       \begin{array}{ll} \textbf{double f2} = -0 \texttt{x1p-1}; \end{array}
30
       double f1 = 0x1p+0;
31
       if (x > 0x1.fe7814e49392fp-1) && (x<1)) {
32
33
           double xx2 = (x-1) * (x-1);
           double ff4 = f4 + f5 * (x-1);
34
           double ff6 = f6 + f7 * (x-1);
35
36
37
           double cr_log_fast_1 = ff4 + xx2 * ff6;
38
39
           /* Add x by -1 result hr,lr */
40
           double hr,lr;
41
           {\tt Add112}({\tt x},\,-1,\,\&{\tt hr},\,\&{\tt lr});
42
43
44
45
           double hr2,1r2;
           /* Multiply (hr,lr) by (hr,lr) result (hr2,lr2) */
46
           Mul222(hr,lr,hr,lr,&hr2,&lr2);
47
48
           double hr4,1r4;
49
50
           /* Multiply (hr2,lr2) by (hr2,lr2) result (hr4,lr4) */
51
           Mu1222(hr2,lr2,hr2,lr2,\&hr4,\&lr4);
52
           double ffh0,ff10;
53
           /* Multiply f1 by (hr,lr) result (ffh0,ffl0) */
54
55
           Mul122(f1,hr,lr,&ffh0,&ffl0);
56
57
           double ffh3,ff13;
           /* Multiply f3 by (hr,lr) result (ffh3,ffl3) */
58
           Mul122(f3,hr,lr,&ffh3,&ffl3);
```

```
60
61
            double ffh2,ff12;
62
            /*Add f2 by (ffh3,ffl3) result (ffh2,ffl2) */
63
            Add122(f2,ffh3,ff13,&ffh2,&ff12);
64
65
            double ffhx2,fflx2;
            /* Multiply (hr2,lr2) by (ffh2,ffl2) result (ffhx2,fflx2) */
66
67
            Mul222(hr2,lr2,ffh2,ffl2,&ffhx2,&fflx2);
68
69
            double ffhx4,fflx4;
70
            /* Multiply (hr4,lr4) by cr_\log_fast_1 result (ffhx4,fflx4) */
 71
            Mul122(cr_log_fast_1,hr4,lr4,&ffhx4,&fflx4);
72
73
            double ffhx0,fflx0;
74
            /* Add (ffh0,ffl0) by (ffhx2,fflx2) result (ffhx0,fflx0) */
            Add222(ffh0, ffl0, ffhx2, fflx2, &ffhx0,&fflx0);
75
76
77
             /* Add (ffhx0,fflx0) by (ffhx4,fflx4) result (h6,l6) */
78
            Add222(ffhx0,fflx0, ffhx4,fflx4,h6,l6);
79
80
            return (*h6+*16);
81
82
83
        }
84
        if ((x> 1) && (x < 0x1.00068db8bac71p+0)){
85
86
87
            double xx2 = (x-1) * (x-1);
            double ff4 = f4 + f5 * (x-1);
88
89
            double ff6 = f6 + f7 * (x-1);
90
91
            double cr_log_fast_1 = ff4 + xx2 * ff6;
92
93
            /* Add x by -1 result hr,lr */
94
95
            double hr,lr;
96
            Add112(x, -1, \&hr, \&lr);
97
98
            double hr2,1r2;
99
            /* Multiply (hr,lr) by (hr,lr) result (hr2,lr2) */
100
101
            Mul222(hr,lr,hr,lr,\&hr2,\&lr2);
102
103
            double hr4,1r4;
            /* Multiply (hr2,lr2) by (hr2,lr2) result (hr4,lr4) */
104
105
            Mul222(hr2,lr2,hr2,lr2,&hr4,&lr4);
106
107
            double ffh0,ff10;
             /* Multiply f1 by (hr,lr) result (ffh0,ffl0) */
108
109
            Mul122(f1,hr,lr,&ffh0,&ffl0);
110
            double ffh3,ff13;
111
            /* Multiply f3 by (hr,lr) result (ffh3,ffl3) */
112
            Mul122(f3,hr,lr,&ffh3,&ffl3);
113
114
115
            double ffh2,ff12;
            /*Add f2 by (ffh3,ffl3) result (ffh2,ffl2) */
116
            Add122(f2,ffh3,ff13,&ffh2,&ff12);
117
118
```

```
119
             double ffhx2,fflx2;
120
             /* Multiply (hr2,lr2) by (ffh2,ffl2) result (ffhx2,fflx2) */
121
             Mul222(hr2,lr2,ffh2,ffl2,&ffhx2,&fflx2);
122
123
             double ffhx4,fflx4;
             /* Multiply (hr4,lr4) by cr log fast 1 result (ffhx4,fflx4) */
124
             \texttt{Mul122}(\texttt{cr\_log\_fast\_1}, \texttt{hr4}, \texttt{lr4}, \texttt{\&ffhx4}, \texttt{\&fflx4});
125
126
127
             double ffhx0,fflx0;
             /* Add (ffh0,ffl0) by (ffhx2,fflx2) result (ffhx0,fflx0) */
128
129
             Add222(ffh0, ff10, ffhx2, ff1x2, &ffhx0,&ff1x0);
130
              /* Add (ffhx0,fflx0) by (ffhx4,fflx4) result (h6,l6) */
131
132
             Add222(ffhx0,fflx0, ffhx4,fflx4,h6,16);
133
134
             return (*h6+*16);
135
136
         }
137
         double m1;
138
139
         /*If x is a subnormal */
         if ((s == 0) \&\& (e == 0) \&\& (m != 0)) {
140
             uint64_t v = m;
141
             e = e - 1023;
142
             v = v*2;
143
             while (v < 0x1000000000000) {
144
145
                 v *= 2;
146
147
148
             m1 = v *0x1p-52;
149
150
151
             u u;
152
             u.x = m1;
             \mathtt{m} = \mathtt{u.i} \ \& \ 0\mathtt{xFFFFFFFFFF};
153
154
         /*If x is normal */
155
         else {
156
157
             m1 = 1 + m*0x1p-52;
158
159
             e = e - 1023;
160
161
         /*We shift by 44 bits for to get the first 8 bits.*/
         {\tt uint64\_t~i} = {\tt m}{>}{>}44;
162
163
164
         alpha_i_m is the double close to 2^k/(2^k+i) such that its log is very close to a double.
165
166
         double alpha_i_m = table_alpha_i_modified[(int)i];
167
168
169
         double hr,lr;
         /*multiply (hm1,lm1) by (halpha_i_m,lalpha_i m)*/
170
         Mul112(m1,alpha_i_m,&hr,&lr);
171
172
         double h,1;
173
         /*add (hr, lr) by -1*/
174
175
         Add122(-1,hr,lr,\&h,\&l);
176
```

```
177
178
179
180
            double hh2 = h * h;
            double ff4 = f4 + f5 * h;
181
182
            double ff6 = f6 + f7 * h;
183
184
185
        double cr_log_fast_1 = ff4 +hh2*ff6;
186
        double 112,hh4,114;
187
188
        /* Multiply (h,l) by (h,l) result (hh2,ll2) */
189
190
        Mul222(h,1,h,1,&hh2,&112);
191
        /* Multiply (hh2,ll2) by (hh2,ll2) result (hh4,ll4) */
192
193
        Mul222(hh2,112,hh2,112,&hh4,&114);
194
195
        double ffh0,ff10;
        /* Multiply f1 by (h,l) result (ffh0,ffl0) */
196
197
        Mul122(f1,h,l,&ffh0,&ffl0);
198
199
        double ffh3,ff13;
200
        /* Multiply f3 by (h,l) result (ffh3,ffl3) */
201
        {\tt Mul122(f3,h,l,\&ffh3,\&ffl3)};\\
202
203
        double ffh2,ff12;
        /*Add f2 by (ffh3,ffl3) result (ffh2,ffl2) */
204
        Add122(f2,ffh3,ff13,&ffh2,&ff12);
205
206
        double ffhx2,fflx2;
207
        /* Multiply (hh2,ll2) by (ffh2,ffl2) result (ffhx2,fflx2) */
208
209
        Mul222(hh2,112,ffh2,ffl2,&ffhx2,&fflx2);
210
211
        double ffhx4,fflx4;
        /* Multiply (hh4,ll4) by cr_log_fast_1 result (ffhx4,fflx4) */
212
213
        Mul122(cr_log_fast_1,hh4,ll4,&ffhx4,&fflx4);
214
215
        double ffhx0,fflx0;
        /* Add (ffh0,ffl0) by (ffhx2,fflx2) result (ffhx0,fflx0) */
216
        Add222(ffh0, ff10, ffhx2, ff1x2, &ffhx0,&ff1x0);
217
218
219
        double h4,14;
        /* Add (ffhx0,fflx0) by (ffhx4,fflx4) result (h4,l4) */
220
        Add222(ffhx0,fflx0, ffhx4,fflx4,&h4,&l4);
221
222
223
        double log_alpha_i_m = table_log_alpha_i_modified[i];
224
225
226
        double h5,15;
227
        /* add (hlog_alpha_i_m, llog_alpha_i_m)*/
        Add122(log_alpha_i_m, h4, 14, &h5, &15);
228
229
230
231
232
233
        double h_{\log 2} = 0x1.62e42fefa38p-1;
        double 1_{\log 2} = 0x1.ef35793c7673p-45;
234
235
        double e_hlog2,e_llog2;
```

crlogaccurate

```
Algorithm 13 Algorithm crlog_{accurate} for x is not to close to 1
Condition: x is a double number and is not to close to 1.
Condition: m_1 and i are results made after calculating x.
Condition: table_{\alpha_i} table_{\log(\alpha_i)} are of tables of Double-Double and are
    calculated with Tang.
Condition: log(2) in double-double: h_{log 2} = 0x1.62e42fefa39efp - 1
    and l_{\log 2} = 0x1.abc9e3b39803fp - 56
Input: m_1 is a double number, i is a Integer number. and table_{\alpha_i} and
    table_{\log(\alpha_i)} are tables of double-double numbers and f a polynomial
    of degree 11 with coefficient f_i calculated with Sollya with 1 \le j \le 11
    and transformed into double-double, G is table of f_i
Condition: 1 \le m_1 \le 2, \ 0 \le i < 256 and size of table_{\alpha_i} = 256 and size of
     table_{\log(\alpha_i)} = 256, \alpha_i = table_{h\alpha_i,l\alpha_i}[i] and \log_{\alpha_i} = table_{h\log(\alpha_i),l\log(\alpha_i)}[i]
Output: \log(x) in double number.
 1: h_r, \ell_r = Mul122(m_1, \alpha_i)
 2: h, \ell = Add122(-1, h_r, \ell_r)
 3: h_1, \ell_1 = Mul222(G[0][0], G[0][1], h, \ell)
 4: for (i = 1, i < 11, i + +)
 5: h_2, \ell_2 = Add222(G[i][0], G[i][1], h_1, \ell_1)
 6: h_1, \ell_1 = Mul222(h_2, \ell_2, h, \ell)
 7: h_5, \ell_5 = Add222(\log_{\alpha_i}, h_1, \ell_1)
 8: e_{hlog2}, e_{\ell log2} = Mul122(e, h_{log2}, \ell_{log2})
 9: h_6, \ell_6 = Add222(e_{hlog2}, e_{\ell log2}, h_5, \ell_5)
10: return (h_6, \ell_6)
```

```
1 # Input : x is a 53-bit floating-point numbers
2 # F is a list with Coefficient of f decomposed into double, double
3 def cr_log_accurate_path(x,F):
4
5
       \# if x<1 and close to 1 return with the function f
6
       if x > RR(0x1.fe73451b9c74fp-1,16) and x < 1:
7
8
           \# we calculated len(F)
9
           lenf = len(F)
10
11
           hr, lr = Add112(-1,x)
12
13
14
           \# F[i] is the coefficient of each monomial :x^{(len(F)-1-i)}
15
           #We are turning F[0] in hf,lf with hf,: main value
```

```
17
            \# lf: error value

hf, lf = RR(F[0][0], 16), RR(F[0][1], 16)

18
19
            \# multiply (hr,lr) by (hf,lf) result (h,l)
20
            h,l = Mul222(hr,lr,hf,lf)
21
22
23
            for i in range(1,lenf):
                #We are turning F[i] in hfi,lfi with hfi: main value
24
25
                 # lfi: error value
26
                hfi,lfi = RR(F[i][0],16),RR(F[i][1],16)
27
28
                 \# add (hfi,lfi) by (h,l) result (h1,l1)
29
                \mathtt{h1,\!11} = \mathtt{Add222}(\mathtt{hfi,\!lfi,\!h,\!l})
30
31
                 \# multiply (hr,lr) by (h1,l1) result (h,l)
                h,l = Mul222(hr,lr,h1,l1)
32
33
34
            return h,1
35
36
37
         \# if x>1 and close to 1 return with function f
38
       if x > 1 and x <= RR('0x1.0178e5916f543p+0',16):
39
40
41
42
            \# we calculated len(F)
43
            lenf = len(F)
44
            hr, lr = Add112(x, -1)
45
46
            \# \ F[i] is the coefficient of each monomial :x^{(len(F)-1-i)}
47
48
            #We are turning F[0] in hf,lf with hf,: main value
49
50
            # lf: error value
            {\tt hf,lf} = {\tt RR}({\tt F[0][0],16}), {\tt RR}({\tt F[0][1],16})
51
52
53
            \# multiply (hr,lr) by (hf,lf) result (h,l)
            h,l = Mul222(hr,lr,hf,lf)
54
55
            for i in range(1,lenf):
56
57
                #We are turning F[i] in hfi,lfi with hfi: main value
                 # lfi: error value
58
59
                hfi,lfi = RR(F[i][0],16),RR(F[i][1],16)
60
                 \# add (hfi,lfi) by (h,l) result (h1,l1)
61
                h1,l1 = Add222(hfi,lfi,h,l)
62
63
64
                 \# multiply (hr,lr) by (h1,l1) result (h,l)
                \mathtt{h,l} = \mathtt{Mul222}(\mathtt{hr,lr,h1,l1})
65
66
            return h,1
67
68
69
        \# s represents the sign
        #e represents the exposant
70
        #m represents fraction
71
72
        (s,m,e)=RR(x).sign_mantissa_exponent()
        \# e = e - 53
73
74
```

```
76
                    \mathtt{e}=\mathtt{e}{+}53
  77
  78
                     \# If x is a subnormal
  79
                    if s==1 and e<=0 and m!=0:
  80
  81
                               \mathtt{v}=\mathtt{m}
  82
                               \mathbf{e} = \mathbf{e}{-}1
  83
                               while v < 2^52:
  84
  85
                                         v*=2
  86
  87
                              m1 = v*2.^(-52)
  88
  89
  90
                     \# {\tt If \ x \ is \ normal}
  91
  92
                     else:
                               \#print(s,e,m)
  93
  94
                              m1 = m*2.^(-52)
                              \mathbf{e}=\mathbf{e}-1
  95
  96
  97
  98
  99
                    binary = RR(m1).str(2)[2:10] \# we recover the 8 bits after the initial 1.
100
101
                     i = int(binary,2) \# i is the 8-bit integer
102
103
                    \verb|halpha_i|, \verb|lalpha_i| = \verb|RR(table_alpha_i[i][0], 16), \verb|RR(table_alpha_i[i][1], 16) \# table
104
                                 computed for all i of alpha_i_m
105
                                                                                                            \# such that log(alpha_i_m) is 71 bits accurate.
                    \verb|hlog_alpha_i|, \verb|llog_alpha_i| = \verb|RR(table_log_alpha_i[i]|[0], 16), \|RR(table_log_alpha_i[i]|[0], 16), \|RR(table_log_alpha_i[i]|[0]
106
                                 [[1],16)
107
                     # multiply (halpha_i_m,lalpha_i_m) by m1
108
                    hr, lr = Mul122(m1, halpha_i, lalpha_i)
109
110
                     \# add hr,lr by (-1) in double,double
111
112
                    h,l = Add122(-1.0,hr,lr)
113
                     # we calculated len(f.list())
114
115
                    lenf = len(F)
116
                     \# F[i] is the coefficient of each monomial :x^{(len(F)-1-i)} in double, double
117
118
                     #We are turning F[0] in hf,lf with hf,: main value
119
120
                     # lf: error value
121
                    hf,lf = RR(F[0][0],16),RR(F[0][1],16)
122
123
                     # multiply (hr,lr) by (hf,lf) result (h,l)
124
                    h1,l1 = Mul222(hf,lf,h,l)
125
                    for i in range(1,lenf):
126
                               #We are turning F[i] in hfi,lfi with hfi: main value
127
128
                               # lfi: error value
129
                               hfi,lfi = RR(F[i][0],16),RR(F[i][1],16)
130
                               \# add (hfi,lfi) by (h1,l1) result (h2,l2)
131
132
                              h2,12 = Add222(hfi,lfi,h1,l1)
```

```
133
            \# multiply (h,1) by (h2,12) result (h1,11)
134
135
            h1,11 = Mul222(h2,12,h,1)
136
137
138
139
        \# Add (h1,l1) by (hlog_alpha_i,llog_alpha_i
140
        h5,15 = Add222(hlog_alpha_i,llog_alpha_i,h1,l1)
141
142
        \# \log(2) = (h_{\log 2, 1_{\log 2}})
143
144
       h_{\log 2} = RR(0x1.62e42fefa39efp-1.16)
        1_{\log 2} = RR(0x1.abc9e3b39803fp-56,16)
145
146
147
        \#e*log(2)
148
        elog2_h,elog2_l = Mul122(RR(e),h_log2,l_log2)
149
150
       hlog_fi, llog_fi = Add222(elog2_h, elog2_l, h5, 15)
151
152
153
        #output: hlog_fi and llog_fi are 53-bit floating-point numbers.
154
155
        # hlog_fi: main value and llog_fi: error value
156
157
        return hlog_fi,llog_fi
  1 static double cr_log_accurate_path(double x, double *h6, double *16){
        int s;
       int e;
 3
       uint64_t m;
 4
 5
        extract (x,\&s,\&e,\&m);
 6
    /* Special case*/
       if (x == 0x1.fb85251a3f26fp-1){
           *h6 = -0x1.1ff9b8e8b38bep-7;
 8
 9
            *16 = -0x1.0b393919c1fa3p-109;
 10
           return (*h6+*16);
 11
       }
 12
        if (x == 0x1.fc65aa1908a66p-1){
 13
            *h6 = -0x1.cecc4ad8d358bp-8;
 14
            *16 = -0x1.65a43e3cf2b61p-107;
 15
 16
            return (*h6+*16);
 17
 18
 19
        /*The coefficients in double, double of the approximation polynomial of degree 11. */
 20
        static const double G[11][2]= {
            \{0 \verb|x1.6d24dd22a92bp-4|, 0 \verb|x1.3e16deaf82f98p-58|\}, \ /* \ \verb|g11| \ decomposed \ into \ double,
 21
            \{-0x1.99889bcf944ecp-4, -0x1.14f5635667f7p-58\}, /* g10 decomposed into
 22
                double, double*
            \{0x1.c71c5b2ab9b1bp-4, -0x1.bba8402b8dab8p-58\}, /* g9 decomposed into double,
 23
                double*/
            \{-0x1.ffffffed56645p-4, 0x1.fa475383e2fep-60\}, /* g8 decomposed into double,
 24
                double*
            \{0x1.249249248d599p-3, 0x1.a5c92da9d350cp-57\}, /* g7 decomposed into double,
 25
                double*.
 26
```

```
double*/
             \{0x1.99999999999ap-3, -0x1.affa1a4be0d04p-57\}, /* g5 decomposed into double,
27
                  double*/
28
              -0x1p-2, 0x1.590555132p-72, /* g4 decomposed into double,double*/
             \{0x1.5555555555555p-2, 0x1.55555540b5114p-56\}, /* g3 decomposed into double,
29
30
             \{-0x1p-1, 0x1.b78p-98\}, /* g2 decomposed into double, double*/
31
             \{0x1p+0, 0x0p+0\} / * g1 decomposed into double, double*/
32
        };
33
34
35
        if ((x > 0x1.fe73451b9c74fp-1) && (x<1)) {
36
            double hr,lr;
             /*Add x by -1 result hr,lr*/
37
38
             Add112(-1,x,&hr,&lr);
39
40
            double hf = G[0][0];
            double lf = G[0][1];
41
42
            double h,1;
43
44
             /*Multiply (hf,lf) by r result (h,l)*/
            Mul222(hr,lr,hf,lf,&h,&l);
45
46
47
            double hfi,lfi,h1,l1;
48
49
            for (int i = 1; i < 11; i++){
                 egin{aligned} \mathbf{hfi} &= \mathbf{G[i][0]}; \\ \mathbf{lfi} &= \mathbf{G[i][1]}; \end{aligned}
50
51
52
53
                 /*Add (hfi,lfi) by (h,l) result (h1,l1)*/
                 Add222(hfi,lfi,h,l,&h1,&l1);
54
55
                 /*multiply (hr,lr) by (h1,l1) result (h,l)*/
56
                 Mul222(hr,lr,h1,l1,&h,&l);
57
58
            }
59
60
             *h6 = h;
61
            *16 = 1;
62
            return (h+1);
63
64
65
        if ((x>1) \&\& (x <= 0x1.0178e5916f543p+0)){
66
67
            double hr,lr;
             /*Add x by -1 result r*/
68
69
             Add112(x,-1,&hr,&lr);
70
71
            double hf = G[0][0];
            double lf = G[0][1];
72
73
74
            double h,1;
             /*Multiply (hf,lf) by r result (h,l)*/
75
            Mul222(hr,lr,hf,lf,\&h,\&l);
76
77
            double hfi,lfi,h1,l1;
78
79
             for (int i = 1; i < 11; i++){
80
                 egin{aligned} \mathtt{hfi} &= \mathtt{G[i][0]}; \\ \mathtt{lfi} &= \mathtt{G[i][1]}; \end{aligned}
81
82
```

```
83
                 /*Add (hfi,lfi) by (h,l) result (h1,l1)*/
 84
 85
                 Add222(hfi,lfi,h,l,&h1,&l1);
 86
                 /*multiply (hr,lr) by (h1,l1) result (h,l)*/
 87
 88
                 Mul222(hr,lr,h1,l1,&h,&l);
 89
 90
             *h6 = h:
 91
 92
             *16 = 1;
 93
 94
             return (h+1);
 95
 96
         double m1;
 97
         /*If x is a subnormal */
 98
         if ((s == 0) \&\& (e == 0) \&\& (m != 0)) {
 99
             uint64_t v = m;
100
             e = e - 1023;
             \mathtt{v}=\mathtt{v*}2;
101
             while (v < 0x1000000000000) {
102
103
                 v *= 2;
                 e--;
104
105
106
            {\tt m1} = {\tt v} *0 {\tt x1p} -52;
107
108
109
110
             \mathtt{u.x}=\mathtt{m1};
            m = u.i \& 0xFFFFFFFFFF;
111
112
         /*If x is normal */
113
        else {
114
            m1 = 1 + m*0x1p-52;
115
116
             e = e - 1023;
117
118
         }
119
         /*We shift by 44 bits for to get the first 8 bits.*/
120
121
        uint64_t i = m >> 44;
122
123
         halpha_i,lalpha_i is the double,double of 2^k/(2^k+i).*/
124
125
        double halpha_i = table_alpha_i[(int)i][0];
        double lalpha_i = table_alpha_i[(int)i][1];
126
127
128
         double hr,lr;
         /*Multiply (halpha_i,lalpha_i_m) by m1 result (hr,lr) */
129
130
        Mul122(m1,halpha_i, lalpha_i,&hr,&lr);
131
132
         double h,1;
         /*Add (hr,lr) by (-1) result (h,l)*/
133
134
         Add122(-1,hr,lr,\&h,\&l);
135
        double hf = G[0][0];
136
        double lf = G[0][1];
137
138
         double h1,11;
139
         /*Multiply (hf,lf) by (h,l) result (h1,l1)*/
140
141
        Mul222(hf,lf,h,l,&h1,&l1);
```

```
142
         double hfi,lfi,h2,l2;
143
144
         for (int i = 1; i < 11; i++){
145
              \mathtt{hfi} = \mathtt{G}[\mathtt{i}][0];
146
              {\tt lfi} = {\tt G[i][1]};
147
148
              /*Add (hfi,lfi) by (h1,l1) result (h2,l2)*/
149
              Add222(hfi,lfi,h1,l1,&h2,&l2);
150
151
              /*multiply (h,l) by (h2,l2) result (h1,l1)*/
152
153
              Mul222(h2,12,h,1,&h1,&l1);
154
         }
/*k=8
155
156
         hlog_alpha_i,llog_alpha_i is the double,double of \log(2^k/(2^k+i)).*/
157
         double hlog_alpha_i = table_log_alpha_i[(int)i][0];
158
         double llog_alpha_i = table_log_alpha_i[(int)i][1];
159
160
161
162
         double h5,15;
         /*Add (h1,l1) by (hlog_alpha_i,llog_alpha_i) result (h5,l5)*/
163
         Add222(hlog_alpha_i,llog_alpha_i,h1,l1,&h5,&l5);
164
165
         \label{eq:double_h_log2} \begin{split} \textbf{double} \ \textbf{h\_log2} &= 0 \texttt{x} 1.62 \texttt{e} 42 \texttt{fefa} \texttt{39efp} {-} 1; \end{split}
166
167
         double 1_{\log 2} = 0x1.abc9e3b39803fp-56;
168
         double e_hlog2,e_llog2;
          /* e*log(2) in (double, double) precision.*/
169
170
         Mul122(e, h_log2, l_log2, &e_hlog2, &e_llog2);
171
172
173
         /*adding (e_hlog2,e_llog2) by (h5,l5) result (h6,l6)*/
174
175
         Add222(e_hlog2, e_llog2, h5, 15, h6,16);
176
         return (*h6+*16);
177
178
179
180 }
```

5

6 8

9 10

11

12

if x > RR(0x1.fba5e353f7ceep-1.16) and x < 1:

we calculated len(F)

hr, lr = Add112(-1.0,x)

lenf = len(F)

crlogadvanced

```
Algorithm 14 Algorithm crlog_{advanced} for x is not to close to 1
  Condition: x is a double number and is not to close to 1.
  Condition: m_1 and i are results made after calculating x.
  Condition: table_{\alpha_i} table_{\log(\alpha_i)} are of tables of triple-Double and are cal-
       culated with Tang.
  Condition: log(2) in triple-double: h_{log 2} = 0x1.62e42fefa39efp - 1,
       lm_{\log 2} = 0x1.abc9e3b39803fp - 56 and l_{\log 2} = 0x1.7b57a079a1934p -
  Input: m_1 is a double number, i is a Integer number. and table_{\alpha_i} and
       table_{\log(\alpha_i)} are tables of triple-double numbers and f a polynomial of
       degree 13 with coefficient f_j calculated with Sollya with 1 \leq j \leq 13
       and transformed into triple-double, U is table of f_i
  Condition: 1 \leq m_1 \leq 2, 0 \leq i < 256 and size of table_{\alpha_i} = 256
       and size of table_{\log(\alpha_i)} = 256, \alpha_i = table_{h\alpha_i, m\alpha_i, l\alpha_i}[i] and \log_{\alpha_i} =
       table_{h \log(\alpha_i), m \log(\alpha_i, l \log(\alpha_i))}[i]
  Output: \log(x) in double number.
   1: h_r, m_r, \ell_r = Mul133(m_1, \alpha_i)
   2: h, m, \ell = Add133(-1, h_r, m_r, \ell_r)
   3: h_1, m_1, \ell_1 = Mul222(U[0][0], U[0][1], U[0][2], h, m, \ell)
   4: for (i = 1, i < 13, i + +)
       h_2, \ell_2 = Add222(U[i][0], U[i][1], U[i][2], h_1, m_1, \ell_1)
       h_1, m_1, \ell_1 = Mul222(h_2, m_2, \ell_2, h, m, \ell)
   7: h_5, m_5, \ell_5 = Add222(\log_{\alpha_i}, h_1, m_1, \ell_1)
   8: e_{hlog2}, e_{mlog2}, e_{\ell log2} = Mul122(e, h_{log2}, m_{log2}, \ell_{log2})
   9: h_6, m_6, \ell_6 = Add222(e_{hlog2}, e_{mlog2}, e_{\ell log2}, h_5, m_5, \ell_5)
  10: m_{\ell} = m_6 + \ell_6 \text{ return } (h_6, m_{\ell})
1 # Input : x is a 53—bit floating—point numbers
2 \# F is a list with Coefficient of f decomposed into triple—double
3 def cr_log_accurate_path_advanced(x,F):
      \# if x>1.17549435*10^(-38) and x<1 \stackrel{	ext{return}}{	ext{return}} with the function f
```

```
14
            \# F[i] is the coefficient of each monomial :x^{(len(F)-1-i)}
15
16
            #We are turning F[0] in hf,lf with hf,: main value
17
            # lf: error value
           hf, mf, lf = RR(F[0][0], 16), RR(F[0][1], 16), RR(F[0][2], 16)
18
19
20
            \# multiply (hr,lr) by (hf,lf) result (h,l)
21
           h,m,l = Mul233(hr,lr,hf,mf,lf)
22
23
           for i in range(1,lenf):
                #We are turning F[i] in hfi,lfi with hfi: main value
24
25
                # lfi: error value
                {\tt hfi,mfi,lfi} = {\tt RR}({\tt F[i][0],16}), {\tt RR}({\tt F[i][1],16}), {\tt RR}({\tt F[i][2],16})
26
27
28
                \# add (hfi,lfi) by (h,l) result (h1,l1)
                h1,m1,l1 = Add333(hfi,mfi,lfi,h,m,l)
29
30
31
                \# multiply (hr,lr) by (h1,l1) result (h,l)
32
                h,m,l = Mul233(hr,lr,h1,m1,l1)
33
34
           return h,m,1
35
36
37
       \# if x> 1.17549435*10^(-38) and x<1 return with the function f
38
39
       if x > 1 and x <= RR(0x1.0178e5916f543p+0.16):
40
41
            # we calculated len(F)
42
43
           lenf = len(F)
44
           hr, lr = Add112(x, -1.0)
45
46
47
           \# F[i] is the coefficient of each monomial :x^{(len(F)-1-i)}
48
49
            \# We are turning F[0] in hf,lf with hf,: main value
50
            \# lf: error value
           hf,mf,lf = RR(F[0][0],16),RR(F[0][1],16),RR(F[0][2],16)
51
52
            \# multiply (hr,lr) by (hf,lf) result (h,l)
53
           h,m,l = Mul233(hr,lr,hf,mf,lf)
54
55
56
           for i in range(1,lenf):
                #We are turning F[i] in hfi,lfi with hfi: main value
57
58
                # lfi: error value
                hfi,mfi,lfi = RR(F[i][0],16),RR(F[i][1],16),RR(F[i][2],16)
59
60
61
                \# add (hfi,lfi) by (h,l) result (h1,l1)
62
                \verb|h1,m1,l1| = \verb|Add333| (\verb|hfi,mfi,lfi,h,m,l)|
63
                \# multiply (hr,lr) by (h1,l1) result (h,l)
64
65
                h,m,1 = Mul233(hr,lr,h1,m1,l1)
66
67
           return h,m,1
68
69
       \#s represents the sign
70
       \#e represents the exposant
71
       \# m represents fraction
72
       (s,m,e) = RR(x).sign_mantissa_exponent()
```

```
73
         \#\mathtt{e}{=}\mathtt{e}{-}53
 74
 75
 76
         \mathbf{e}=\mathbf{e}{+}53
 77
 78
         \# If x is a subnormal
         if s==1 and e<=0 and m!=0:
 79
 80
 81
             \mathtt{v}=\mathtt{m}
 82
              \mathbf{e}=\mathbf{e}{-}1
 83
 84
              while v < 2^52:
                  v*=2
 85
 86
 87
                  e=e-1
             m1 = v*2.^(-52)
 88
 89
 90
        \#If x is normal
 91
         else:
 92
             \#print(s,e,m)
 93
             m1 = m*2.^(-52)
             e = e-1
 94
 95
         binary = RR(m1).str(2)[2:10] # we recover the 8 bits after the initial 1.
 96
 97
         i = int(binary,2) \# i \text{ is the } 8-bit \text{ integer}
 98
 99
100
         \mathtt{halpha\_i,malpha\_i,lalpha\_i} = \mathtt{RR}(\mathtt{table\_alpha\_i\_triple[i][0],} 16),\mathtt{RR}(\mathtt{lalpha\_i,malpha\_i,lalpha\_i})
               \verb|table_alpha_i_triple[i][1],16), \verb|RR(table_alpha_i_triple[i][2],16)| \# table|
              computed for all i of alpha_i_m
                                                \# such that log(alpha_i_m) is 71 bits accurate.
101
102
         \verb|hlog_alpha_i,mlog_alpha_i,1log_alpha_i = RR(table_log_alpha_i\_triple[i][0],16),RR(table_log_alpha_i,mlog_alpha_i)|
              table_log_alpha_i_triple[i][1],16),RR(table_log_alpha_i_triple[i][2],16)
103
         # multiply (halpha_i_m,lalpha_i_m) by m1
104
         hr,mr,lr = Mul133(m1,halpha_i,malpha_i,lalpha_i)
105
106
107
         \# add hr,mr,lr by (-1) in triple-double
108
         h,m,l = Add133(-1.0,hr,mr,lr)
109
         # we calculated len(f.list())
110
111
         lenf = len(F)
112
         \# F[i] is the coefficient of each monomial :x^{(en(F)-1-i)} in double, double
113
114
         #We are turning F[0] in hf,lf with hf,: main value
115
116
         # lf: error value
117
         hf,mf,lf = RR(F[0][0],16),RR(F[0][1],16),RR(F[0][2],16)
118
119
         # multiply (hr,lr) by (hf,lf) result (h,l)
120
         h1,m1,l1 = Mul333(h,m,l,hf,mf,lf)
121
199
         for i in range(1,lenf):
123
              #We are turning F[i] in hfi,lfi with hfi: main value
124
              # lfi: error value
             hfi,mfi,lfi = RR(F[i][0],16),RR(F[i][1],16),RR(F[i][2],16)
125
126
127
              \# add (hfi,lfi) by (h1,l1) result (h2,l2)
128
              h2,m2,12 = Add333(hfi,mfi,lfi,h1,m1,l1)
```

```
129
            \# multiply (h,1) by (h2,12) result (h1,11)
130
131
            h1,m1,l1 = Mul333(h,m,l,h2,m2,l2)
132
133
134
135
        \# Add (h1,l1) by (hlog_alpha_i,llog_alpha_i
136
        h5,m5,15 = Add333(hlog_alpha_i,mlog_alpha_i,llog_alpha_i,h1,m1,l1)
137
138
        \# \log(2) = (h_{\log 2}, l_{\log 2})
139
140
        h_{\log 2} = RR(0x1.62e42fefa39efp-1.16)
        m_log2 = RR(0x1.abc9e3b39803fp-56,16)
141
142
        1_{\log 2} = RR(0x1.7b57a079a1934p-111,16)
143
144
        \#e*log(2)
145
        \verb|elog2_h,elog2_m,elog2_l| = \verb|Mul133(RR(e),h_log2,m_log2,l_log2)|
146
147
148
        hlog_fi,mlog_fi,llog_fi = Add333(elog2_h,elog2_m,elog2_1,h5,m5,15)
149
150
151
        #output: hlog_fi and llog_fi are 53-bit floating-point numbers.
152
        # hlog_fi: main value and llog_fi: error value
153
154
        return hlog_fi, mlog_fi, llog_fi
 1 static double cr_log_accurate_advanced_path(double x){
 2
        int s;
 3
        int e;
 4
        uint64_t m;
 5
        double h6,m1;
 6
        /*special cases*/
 7
        if (x == 0x1.c7e1077f9aec2p-1) {
 8
            {\tt h6} = -0{\tt x1.db88b1160524bp}{-4};
 9
            10
            return (h6+ml);
11
        if (x == 0x1.f8db13b0e98a3p-1) {
12
           h6 = -0x1.cc7365166e597p-7;
13
            ml = 0x1.fffffffffffccp-61;
 14
            return (h6+ml);
15
16
        if (x == 0x1.acb8cf13bc769p-2) 
17
            h6 = -0x1.bdc7955d1482cp-1;
18
 19
            {\tt ml} = 0 {\tt x1.6c68c50e2ff8p-} 105;
20
            return (h6+ml);
 21
        if (x == 0x1.9309142b73ea6p-1) {
22
 23
            h6 = -0x1.ea16274b0109bp-3;
            {\tt ml} = 0 {\tt x1.705633d4614a5p-106};
24
            return (h6+ml);
25
26
        if (x == 0x1.98a04e0833091p-1) {
27
28
           h6 = -0x1.cddf723d3e52fp-3;
           ml = 0x1.a0636ba8c84a5p-106;
29
30
            return (h6+ml);
```

```
31
      if (x == 0x1.baded30cbf1c4p-1) {
32
33
          h6 = -0x1.290ea09e36479p-3;
34
          ml = 0x1.3339fcb8b1c92p-111;
35
           return (h6+ml);
36
37
      if (x == 0x1.c7f14af0a08ebp-1) {
38
          h6 = -0x1.daf693d64fadap-4;
39
          ml = 0x1.4f60080a4d88dp-109;
          return (h6+ml);
40
      }
41
42
      if (x == 0x1.f3c35328f1d5dp-1) {
43
          h6 = -0x1.8c56ff5326197p-6;
44
45
          ml = 0x1.24265dc55252p-105;
46
           return (h6+ml);
47
48
      if (x == 0x1.f8156947924c8p-1) {
49
          h6 = -0x1.fe9ad6761218dp-7;
50
51
          ml = 0x1.39b6b04fc1d53p-105;
52
           return (h6+ml);
53
      }
54
55
      extract (x,&s,&e, &m);
56
57
      /*The coefficients in triple—double numbers of the approximation polynomial of degree
58
           14. */
59
      static const double U[14][3]= {
          \{-0 \verb|x1.1d35b84c5cc56p-4|, 0 \verb|x1.780337dcf5705p-58|, -0 \verb|x1.b91450d691e2p-114|\}, /* \\
60
              u14 decomposed into triple—double*
          \{0x1.3afccbd40d20fp-4, 0x1.9c4dd4dab116p-58, 0x1.3a1e112af7efp-112\}, /* u13
61
              decomposed into triple-double*/
          \{-0x1.55552b4de4c76p-4, -0x1.50aa78415c4bfp-59, 0x1.2fafbf95a1a8p-115\}, /*
62
              u12 decomposed into triple-double*/
          \{0x1.745d171352e2dp-4, -0x1.b9a4b0517c2c6p-58, -0x1.6128b0932efc8p-112\}, /*
63
               u11 decomposed into triple—double*,
64
          \{-0x1.999999996fed4p-4, 0x1.8ec6a7a38be73p-58, -0x1.85507f01f5868p-112\}, /*
              u10 decomposed into triple—double*,
          \{0x1.c71c71c71c59ap-4, 0x1.1e1945381ebebp-58, -0x1.b30f3065fce38p-112\}, /*
65
              u9 decomposed into triple—double*,
          \{-0x1.ffffffffffffp-4, -0x1.89d824c2891b7p-58, 0x1.b4f656bbfa58p-115\}, /*
66
              u8 decomposed into triple—double*/
          \{0x1.2492492492492p-3, 0x1.2438b8708dbd8p-57, 0x1.e00436fa1c45p-113\}, /* u7\}
67
              decomposed into triple-double*/
          \{-0x1.5555555555555555p-3, -0x1.5555440ddf7efp-57, 0x1.0fce30dd1d9p-116\}, /*u6\}
68
               decomposed into triple-double*
          69
              decomposed into triple-double*/
          \{-0x1p-2, 0x1.4b03dac3ad4d9p-100, 0x1.78p-156\}, /* u4 decomposed into triple-
70
          71
              decomposed into triple—double*/
           -0x1p-1, 0x1.9745a799cp-127, 0x0p+0\}, /* u2 decomposed into triple-double*/
72
          \{0x1p+0, -0x1.09f14p-143, 0x0p+0\}, /* u1 decomposed into triple—double*/
73
74
      };
75
     if ((x>0x1.fc65aa1908a66p-1) \&\& (x<1)) {
```

```
77
             double hr,lr ;
              /*Add x by -1 result hr,lr*/
 78
 79
             Add112(-1.0,x,\&hr,\&lr);
 80
            double hf = U[0][0];
 81
             double mf = U[0][1];
 82
             double lf = U[0][2];
 83
 84
 85
             double h,m,1;
 86
              /*Multiply (hf,mf,lf) by r result (h,m,l)*/
             {\tt Mul233(hr,lr,hf,mf,lf,\&h,\&m,\&l)};\\
 87
 88
              /*printf("i = 0, h = \%la, m = \%la, l = \%la \n",h,m,l);*/
 89
             double hfi,mfi,lfi,h1,m1,l1;
 90
 91
             for (int i = 1; i < 14; i++){
                  \mathtt{hfi} = \mathtt{U[i][0]};
 92
 93
                  mfi = U[i][1];
                  \texttt{lfi} = \texttt{U[i][2]};
 94
                  /*Add (hfi,mfi,lfi) by (h,m,l) result (h1,m1,l1)*/
 95
                  Add333(hfi,mfi,lfi,h,m,1,&h1,&m1,&l1);
 96
 97
                  /*multiply (hr,lr) by (h1,m1,l1) result (h,m,l)*/
 98
                  {\tt Mul233(hr,lr,h1,m1,l1,\&h,\&m,\&l)};\\
99
100
101
            }
102
103
            return h+m;
         }
104
105
106
         if ((x>1) \&\& (x <= 0x1.0178e5916f543p+0)) {
107
             double hr,lr;
              /*Add x by -1 result hr, lr*/
108
          Add112(x,-1.0,\&hr,\&lr);
109
110
             double hf = U[0][0];
111
             double mf = U[0][1];
112
113
             double lf = U[0][2];
114
115
             double h,m,1;
              /*Multiply (hf,mf,lf) by r result (h,m,l)*/
116
             Mul233(hr,lr,hf,mf,lf,\&h,\&m,\&l);
117
118
             double hfi,mfi,lfi,h1,m1,l1;
119
120
             for (int i = 1; i < 14; i++){
121
                  \mathtt{hfi} = \mathtt{U}[\mathtt{i}][0];
122
                  \mathtt{mfi} = \mathtt{U[i][1]};
123
124
                  {\tt lfi} = {\tt U[i][2]};
                  /*Add (hfi,mfi,lfi) by (h,m,l) result (h1,m1,l1)*/
125
126
                  Add333(hfi,mfi,lfi,h,m,1,&h1,&m1,&l1);
127
                  /*multiply (hr,lr) by (h1,m1,l1) result (h,m,l)*/
128
                  Mul233(hr,lr,h1,m1,l1,&h,&m,&l);
129
130
             }
131
132
133
134
             return h+m;
135
         }
```

```
136
        double m1;
         /*If x is a subnormal */
137
138
        if ((s == 0) \&\& (e == 0) \&\& (m != 0)) {
139
             uint64_t v = m;
            e = e - 1023;
140
             v = v*2;
141
             while (v < 0x100000000000) {
142
143
                v *= 2;
144
                e--:
145
146
147
            m1 = v *0x1p-52;
148
149
            u u;
150
            \mathtt{u.x} = \mathtt{m1};
            151
152
153
         /*If x is normal */
        else {
154
155
            m1 = 1 + m*0x1p-52;
156
             e = e - 1023;
157
158
159
         /*We shift by 44 bits for to get the first 8 bits.*/
160
161
        uint64_t i = m >> 44;
162
163
         /*k=8
164
165
        halpha i,lalpha i is the double,double of 2<sup>k</sup>/(2<sup>k</sup>+i).*/
        double halpha_i = table_alpha_i_triple[(int)i][0];
166
        double malpha_i = table_alpha_i_triple[(int)i][1];
167
168
        double lalpha_i = table_alpha_i_triple[(int)i][2];
169
          double hr,mr,lr;
170
         /*Multiply (halpha i, malpha i, lalpha i m) by m1 result (hr, mr, lr) */
171
172
        Mul133(m1,halpha_i, malpha_i, lalpha_i,&hr,&mr,&lr);
173
174
        double h,mprime,1;
         /*Add (hr,lr) by (-1) result (h,m,l)*/
175
176
        Add133(-1,hr,mr,lr,&h,&mprime,&1);
177
        double hf = U[0][0];
178
        double mf = U[0][1];
179
        double lf = U[0][2];
180
181
        double h1, m1prime, l1;
         /*Multiply (hf,mf,lf) by (h,m,l) result (h1,m1,l1)*/
182
183
        Mul333(h,mprime,l,hf,mf,lf,&h1,&m1prime,&l1);
184
185
        double hfi,mfi,lfi,h2,m2,l2;
186
187
        for (int i = 1; i < 14; i++){
            \mathtt{hfi} = \mathtt{U[i][0]};
188
            mfi = U[i][1];
189
             lfi = U[i][2];
190
             /*Add (hfi,mfi,lfi) by (h1,m1,l1) result (h2,m2,l2)*/
191
192
             Add333(hfi,mfi,lfi,h1,m1prime,l1,&h2,&m2,&l2);
193
194
             /*multiply (h,m,l) by (h2,m2,l2) result (h1,m1,l1)*/
```

```
195
              Mul333(h,mprime,1,h2,m2,12,&h1,&m1prime,&l1);
196
197
198
          /*k=8
         hlog alpha i,llog alpha i is the double, double of \log(2^k/(2^k+i)).*
199
         double hlog_alpha_i = table_log_alpha_i_triple[(int)i][0];
200
201
         \label{eq:double_mlog_alpha_i} \begin{aligned} & \text{double mlog_alpha_i} = \text{table\_log\_alpha\_i\_triple}[(\text{int})\text{i}][1]; \end{aligned}
202
         double llog_alpha_i = table_log_alpha_i_triple[(int)i][2];
203
204
         double h5,m5,15;
205
206
         /*Add (h1,m1,l1) by (hlog alpha i,mlog alpha i,llog alpha i) result (h5,m5,l5)*/
         {\tt Add333(hlog\_alpha\_i,mlog\_alpha\_i,llog\_alpha\_i,h1,m1prime,l1,\&h5,\&m5,\&15)};\\
207
208
209
         double h_log2 = 0x1.62e42fefa39efp-1;
         double m_log2 = 0x1.abc9e3b39803fp-56;
210
211
         double 1_{\log 2} = 0x1.7b57a079a1934p-111;
212
         double e_hlog2,e_mlog2,e_llog2;
213
         /* e*log(2) in triple—double precision.*/
         {\tt Mul133(e, h\_log2, m\_log2, l\_log2, \&e\_hlog2, \&e\_mlog2, \&e\_llog2);}
214
215
         double m6,16;
216
217
218
         /*adding (e hlog2,e mlog2,e llog2) by (h5,m5,l5) result (h6,m6,l6)*/
219
         Add333(e_hlog2, e_mlog2, e_llog2, h5, m5, 15, &h6, &m6, &l6);
220
221
222
         ml = m6+16;
223
224
225
226
         \frac{\text{double m62}}{\text{double m62}} = 2*\text{m6};
227
        if ( (fabs(m6) == ulp_0_5(h6)) \&\& m6>0){
228
              if (16 > 0){
229
230
                  ml = nextafter(m6,m62);
231
232
              if (16<0){
233
                   ml = nextafter(m6,0);
234
235
          \begin{tabular}{ll} \bf if ( (fabs(m6) == ulp_0_5(h6)) \&\& (m6{<}0)) \{ \end{tabular} 
236
237
              if (16<0){
238
                  ml = nextafter(m6,m62);
239
              if (16>0){
240
241
                  ml = nextafter(m6,0);
242
243
244
         return h6+m1;
245 }
```

log

```
Algorithm 15 Algorithm log

Condition: x is a double number.

Condition: err_{fast} = 0x1.6b6b11ea279ecp - 59 and err_{accurate} = 0x1.810c9a86fc45ep - 99

Output: \log(x) in double number.

1: y, h_6, l_6 = cr \log_{fast}(x)

2: left = h6 + FMA(-errfast, h6, l6)

3: right = h6 + FMA(+errfast, h6, l6)

4: ifright == left

5: return y

6: y, h_6, l_6 = cr \log_{accurate}(x)

7: right = h6 + FMA(+errfast, h6, l6)

8: ifright == left

9: return y

10: y = cr \log_{advanced}(x) return y
```

```
1 /* x is double*/
 2 double cr_log(double x){
       double h6 = 0;
 4
       double 16 = 0;
       double right,left;
 5
       double y;
 6
       double errfast = 0x1.6b6b11ea279ecp-59;
 7
       /* We started with the variable errfast = 0x1.01c7d6c404f05p-63 according to
9
       document (CR-LIBM A library of correctly rounded elementary functions in double-
            precision: 2006).
10
       we had errors in all 4 rounding modes.
       We increased the power by 1 until we found only one error in the RNDN rounding mode
11
12
       and no other errors in the other rounding modes. After, we had err = 0x1p-64.
       We searched by dichotomy after the dot. If we had errors on the equality between
            \operatorname{cr} \log \operatorname{fast}(x) and \operatorname{mpfr} \log(x)
        we increase either we decrease until we have 2 consecutive digits.
14
15
       We started again until we optimize errfast and also erracc.
16
17
18
       y = cr_log_fast_path(x,&h6,&l6);
       left = h6 + \__builtin_fma (-errfast, h6, 16);
       right = h6 + __builtin_fma (+errfast, h6, 16);
20
21
22
       if (right== left){
23
24
           return y;
25
       double erracc = 0x1.810c9a86fc45ep-99;
26
       /* We started with the variable errace = 0x1p-102, Result obtained by calculating with
```

```
We do as we did with errfast*/ \label{eq:cr_log_accurate_path} \begin{split} y &= \texttt{cr_log_accurate_path}(x,\&\texttt{h6},\&\texttt{16}); \\ \texttt{left} &= \texttt{h6} + \_\texttt{builtin_fma} \; (-\texttt{erracc},\,\texttt{h6},\,\texttt{16}); \\ \texttt{right} &= \texttt{h6} + \_\texttt{builtin_fma} \; (+\texttt{erracc},\,\texttt{h6},\,\texttt{16}); \end{split}
29
30
31
32
33
                     \quad \quad \textbf{if} \; (\texttt{right} {==} \; \texttt{left}) \{
34
35
                                 return y;
36
37
38
                     {\tt y = cr\_log\_accurate\_advanced\_path(x);}
39
                     return y;
40
41 /* log(x) in double*/
42 }
```

Bibliography

- [1] Mohammed Belaid. Résolution de contraintes sur les flottants dédiée à la vérification de programmes. PhD thesis, Université Nice Sophia Antipolis, 2013.
- [2] Sylvie Boldo and Marc Daumas. Properties of the subtraction valid for any floating point system. In 7th International Workshop on Formal Methods for Industrial Critical Systems, pages 137–149, 2002.
- [3] Sylvain Chevillard, Christoph Lauter, and Mioara Joldes. Users' manual for the sollya tool. 2013.
- [4] Catherine Daramy-Loirat, David Defour, Florent De Dinechin, Matthieu Gallet, Nicolas Gast, Christoph Quirin Lauter, and Jean-Michel Muller. Cr-libm a library of correctly rounded elementary functions in double-precision. Technical report, Technical report, LIP, ENS Lyon, 2006. URL: https://hal-ens-lyon.archives..., 2009.
- [5] Florent De Dinechin, Christoph Lauter, and Jean-Michel Muller. Fast and correctly rounded logarithms in double-precision. *RAIRO-Theoretical Informatics and Applications*, 41(1):85–102, 2007.
- [6] Theodorus Jozef Dekker. A floating-point technique for extending the available precision. *Numerische Mathematik*, 18(3):224–242, 1971.
- [7] Nicolas Fabiano, Jean-Michel Muller, and Joris Picot. Algorithms for triple-word arithmetic. *IEEE Transactions on Computers*, 68(11):1573–1583, 2019.
- [8] Shmuel Gal. Computing elementary functions: A new approach for achieving high accuracy and good performance. In *Accurate scientific computations*, pages 1–16. Springer, 1986.
- [9] Stef Graillat and Fabienne Jézéquel. Tight interval inclusions with compensated algorithms. *IEEE Transactions on Computers*, 69(12):1774–1783, 2019.

124 BIBLIOGRAPHY

[10] Vincenzo Innocente and Paul Zimmermann. Accuracy of mathematical functions in single, double, double extended, and quadruple precision. 2021.

- [11] Claude-Pierre Jeannerod, Jean-Michel Muller, and Paul Zimmermann. On various ways to split a floating-point number. In 2018 IEEE 25th Symposium on Computer Arithmetic (ARITH), pages 53–60. IEEE, 2018.
- [12] Christoph Quirin Lauter. Basic building blocks for a triple-double intermediate format. PhD thesis, INRIA, LIP, 2005.
- [13] Julien Le Maire, Nicolas Brunie, Florent De Dinechin, and Jean-Michel Muller. Computing floating-point logarithms with fixed-point operations. In 2016 IEEE 23nd Symposium on Computer Arithmetic (ARITH), pages 156–163. IEEE, 2016.
- [14] Vincent Lefèvre and Paul Zimmermann. Arithmétique flottante. PhD thesis, INRIA, 2004.
- [15] Nicolas Louvet. Algorithmes compensés en arithmétique flottante: précision, validation, performances. PhD thesis, Université de Perpignan Via Domitia, 2007.
- [16] Peter Markstein. The new ieee-754 standard for floating point arithmetic. In *Dagstuhl Seminar Proceedings*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2008.
- [17] Jean-Michel Muller. Elementary functions.
- [18] Jean-Michel Muller. On the definition of ulp (x). PhD thesis, INRIA, LIP, 2005.
- [19] Jean-Michel Muller, Nicolas Brisebarre, Florent de Dinechin, Claude-Pierre Jeannerod, Vincent Lefèvre, Guillaume Melquiond, Nathalie Revol, Damien Stehlé, and Serge Torres. Handbook of floating-point arithmetic, 2010.
- [20] Deevashwer Rathee, Anwesh Bhattacharya, Rahul Sharma, Divya Gupta, Nishanth Chandran, and Aseem Rastogi. Secfloat: Accurate floating-point meets secure 2-party computation (full version). Technical report, Cryptology ePrint Archive, Report 2022/322, 2022, https://ia.cr/2022/322.

BIBLIOGRAPHY 125

[21] Wikipedia contributors. Ieee 754 — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=IEEE_754&oldid=1096894205, 2022. [Online; accessed 1-September-2022].