



HAL
open science

When $(3x/3)$ and $3(x/3)$ are not equal to x

Frédéric Goualard

► **To cite this version:**

| Frédéric Goualard. When $(3x/3)$ and $3(x/3)$ are not equal to x . 2016. hal-01451457

HAL Id: hal-01451457

<https://hal.science/hal-01451457v1>

Preprint submitted on 1 Feb 2017

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

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

When $(3x)/3$ and $3(x/3)$
are not equal to x

Rapport de recherche
Research Report

RR n° xx.xx

Frédéric GOUALARD

Frédéric GOULARD

When $(3x)/3$ and $3(x/3)$ are not equal to x

Les rapports de recherche du Laboratoire des Sciences du Numérique de Nantes sont disponibles au format Adobe® PDF sur <http://ls2n.fr/>.

Research Reports from the Laboratoire des Sciences du Numérique de Nantes are available at <http://ls2n.fr/>.

© February 2017 by Frédéric GOULARD.

When $(3x)/3$ and $3(x/3)$ are not equal to x

Frédéric GOUALARD

Abstract

Rounding Error Analysis is routinely used to compute a worst-case error bound on the result of algorithms that use floating-point arithmetic. However, for some applications (e.g., when it is necessary to prove some inclusion of the result in a domain), the knowledge of both an upper-bound of the magnitude of the error and of its sign is paramount. Using standard rounding error analysis together with a simple systematic approach, we compute such information for the expressions $3x$, $x/3$, $3(x/3)$, and $3x/3$, which can be used, e.g., in the proof of interval arithmetic operators.

CCS Concepts: Mathematics of computing→Interval arithmetic, *Theory of computation*→*Rounding techniques*

Additional Key Words and Phrases: floating-point arithmetic, error analysis, roundoff, interval arithmetic

Contents

1	Introduction	1
2	Binary floating-point numbers	1
3	Multiplying x by three	2
3.1	The normal case	2
3.2	The denormal case	3
4	Dividing x by three	6
4.1	The normal case	6
4.2	The denormal case	8
5	Composing the multiplication and division	13
5.1	Multiplication and division with no underflow	13
5.2	Multiplication and division with underflow	13
5.2.1	The case of $\text{fl}\langle\langle 3x \rangle / 3 \rangle$	14
5.2.2	The case of $\text{fl}\langle 3 \langle x / 3 \rangle \rangle$	15
5.3	Putting it all together	16
6	Conclusion	18
A	Multiplication by 3 involving no denormal number	21
B	Multiplication by 3 involving denormal numbers	22
C	Division by 3 involving no denormal number	23
D	Division by 3 involving denormal numbers	28

1 Introduction

The implementation of Interval Arithmetic operators requires the evaluation of floating-point expressions, which can be marred by rounding errors that jeopardize the very properties that Interval Arithmetic is supposed to ensure [4]. *Rounding error analysis* [6] can be used to determine in a systematic way the magnitude of these errors. It is usually, however, not concerned with their sign, a fatal flaw whenever the question of the inclusion of a value into an interval is to be ascertained.

Some clever use of the binary floating-point number format properties as guaranteed by the IEEE 754 standard [7] allows in special cases to work out the conditions under which rounding errors cancel out in such a way that the computed result is exactly equal to the true result. Several examples of such an approach can be found in Goldberg’s survey [3] and Muller *et al.*’s book [9], to name a few sources. On the other hand, few works concern themselves with determining both the magnitude *and* the sign of rounding errors in computing arithmetic expressions.

In this paper, we investigate the case of the expressions $3x/3$ and $3(x/3)$ for x a binary floating-point number. Using a mixture of classical Rounding Analysis and systematic study of the elementary operations “ $3x$ ” and “ $x/3$ ”, we determinate for each possible value of x the sign and the magnitude of the error in computing $3x/3$ and $3(x/3)$. This study is used in a forthcoming article [5] to prove inclusion properties of some interval trisection operators. We expect it to be useful in the study of the many other algorithms that rely on a multiplication or a division by 3 (see, e.g., de Dinechin *et al.*’s work [1]).

2 Binary floating-point numbers

In this study, we only concern ourselves with binary floating-point numbers as defined by the 1985 IEEE 754 standard [7], by far the most widespread standard in use today for floating-point numbers. This paper assumes a prior knowledge of the basics of the IEEE 754 standard even though the most important points are summarized hereafter as needed for the sake of completeness.

We consider binary floating-point numbers from a set \mathbb{F} represented with a significand of size p bits:

$$x \in \mathbb{F}, \quad x = \pm b_0.b_1 \dots b_{p-1} \times 2^e$$

with an exponent e in the range $[E_{\min}, E_{\max}]$.

For the purpose of the systematic study of the elementary operations, we assume $p \geq 5$, not a strong constraint considering that the smallest format presented in the IEEE 754 standard has $p = 24$. We will not take into account the possibility of overflow when performing the multiplication by 3; on the other hand, the possibility of underflow will be fully adressed, albeit separately from the normal cases.

When the base b used to represent a number x is not unambiguously drawn from the context, we will indicate it as such: x_b .

All results are supposed to be rounded to nearest-even, the usual default rounding strategy. Given r a real number, we note $\text{fl}(r)$ the floating-point value corresponding to r rounded to nearest-even. Given x , a floating-point number, we note x^- the greatest floating-point number smaller than x , and x^+ the smallest floating-point number greater than x . For an expression of the form $(x \diamond y) \square z$ —with $(x, y, z) \in \mathbb{F}^3$ and (\diamond, \square) arithmetic operators—, we use the shorthand $\text{fl}\langle(x \diamond y) \square z\rangle^1$ to mean $\text{fl}(\text{fl}(x \diamond y) \square z)$. Lastly, we will extensively use in this paper the fact that rounding is order-preserving (*rounding monotonicity* [6, p. 38]):

$$\forall(r_1, r_2) \in \mathbb{R}^2 : r_1 \geq r_2 \implies \text{fl}(r_1) \geq \text{fl}(r_2) \quad (1)$$

Note that if there is a strict inequality between r_1 and r_2 (i.e., $r_1 > r_2$), we still have $\text{fl}(r_1) \geq \text{fl}(r_2)$.

With *correctly rounded* floating-point operators, as specified by the IEEE 754 standard, there is a simple relationship between a computed value and the real result [6]:

$$\text{fl}(x \diamond y) = (x \diamond y)(1 + \delta) + \eta, \quad \begin{cases} (x, y) \in \mathbb{F}^2, \\ \delta \eta = 0, \\ |\delta| \leq u, \\ |\eta| \leq \mu/2 \\ \diamond \in \{+, -, \times, \div\} \end{cases} \quad (2)$$

¹Note the parentheses replaced by chevrons.

where u is the *unit roundoff* equal to 2^{-P} and μ is the *smallest positive subnormal number* equal to $u2^{E_{\min}+1}$. The bound on δ may be somewhat refined depending on the operator; in particular, we will make use of the following improved bounds [8]:

$$\begin{cases} \text{fl}(xy) = xy(1 + \delta) & , \quad |\delta| \leq \frac{u}{1+u} \\ \text{fl}\left(\frac{x}{y}\right) = \frac{x}{y}(1 + \delta) & , \quad |\delta| \leq u - 2u^2 \end{cases} , \quad (x, y) \in \mathbb{F}^2, \quad (3)$$

provided no underflow or overflow occurs.

3 Multiplying x by three

Barring overflow, the multiplication of two binary floating-point numbers:

$$x = (-1)^{s_x} m_x \times 2^{e_x}$$

and

$$y = (-1)^{s_y} m_y \times 2^{e_y}$$

is given by [9]:

$$xy = (-1)^{s_x \oplus s_y} m_x m_y \times 2^{e_x + e_y} \quad (4)$$

The multiplication of a floating-point number x by 3 (i.e., 1.1×2^1 in normal floating-point binary form) is particularly easy to perform thanks to the simplicity of the representation of 3. We will consider first the case for which no underflow occurs; denormal numbers will be considered next.

3.1 The normal case

In this section, x is a normal number of the form $x = (-1)^{s_x} 1.b_1 \dots b_{p-1} \times 2^{e_x}$. From Eq. (4), we see that we may consider $x > 0$ only, as the results thus obtained need only be mirrored to obtain their counterpart when $x < 0$.

Depending on x , there are two possibilities: either $3x$ has a $p+1$ bits significand (Eq. (5a)), or it has a $p+2$ bits significand (Eq. (5b)).

$$\begin{array}{r} \times \quad \begin{array}{cccccccc} & 1. & b_1 & b_2 & \cdots & b_{p-3} & b_{p-2} & b_{p-1} & \times 2^{e_x} \\ & 1. & 1 & & & & & & \times 2^1 \end{array} \\ \hline + \quad \begin{array}{cccccccc} & 1. & b_1 & b_2 & \cdots & b_{p-3} & b_{p-2} & b_{p-1} & \\ & 1. & b'_1 & b'_2 & \cdots & b'_{p-3} & b'_{p-2} & b'_p & \times 2^{e_x+1} \end{array} \end{array} \quad (5a)$$

$$\begin{array}{r} \times \quad \begin{array}{cccccccc} & 1. & b_1 & b_2 & \cdots & b_{p-3} & b_{p-2} & b_{p-1} & \times 2^{e_x} \\ & 1. & 1 & & & & & & \times 2^1 \end{array} \\ \hline + \quad \begin{array}{cccccccc} & 1. & b_1 & b_2 & \cdots & b_{p-3} & b_{p-2} & b_{p-1} & \\ & 1. & b'_1 & b'_2 & \cdots & b'_{p-3} & b'_{p-2} & b'_p & b'_{p+1} & \times 2^{e_x+2} \end{array} \end{array} \quad (5b)$$

From Eqs. (5a) and (5b), we can deduce conditions on x for $y = 3x$ to have a $p+1$ bits significand; similarly, we can deduce the value of b'_1 depending on the number of bits of $3x$:

Lemma 1. *Given $x = 1.b_1 \dots b_{p-1} \times 2^{e_x}$ a normal binary floating-point number. If $y = 3x$ has a $p+1$ bits significand, then $x = 1.0b_2 \dots b_{p-1} \times 2^{e_x}$ and $y = 1.1b'_2 \dots b'_{p-1}b'_p \times 2^{e_x+1}$*

Table 1: Result and sign of the error for the multiplication by 3 of $x = 1.0b_2 \dots b_{p-1} \times 2^{e_x}$

\mathbf{x}	$\text{fl}(3\mathbf{x})$	Error
$1.0b_2 \dots b_{p-4}000$	$1.1b'_2 \dots b'_{p-4}000$	EQ
	$1.1b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}010$	EQ
	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}110$	EQ
$1.0b_2 \dots b_{p-4}001$	$1.1b'_2 \dots b'_{p-4}010$	GT
	$1.1b'_2 \dots b'_{p-4}110$	GT
	$1.0b'_2 \dots b'_{p-4}001$	GT
	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.0b'_2 \dots b'_{p-4}101$	GT
	$1.0b'_2 \dots b'_{p-4}111$	GT
$1.0b_2 \dots b_{p-4}010$	$1.1b'_2 \dots b'_{p-4}011$	EQ
	$1.1b'_2 \dots b'_{p-4}111$	EQ
	$1.0b'_2 \dots b'_{p-4}000$	GT
	$1.0b'_2 \dots b'_{p-4}010$	GT
	$1.0b'_2 \dots b'_{p-4}100$	GT
	$1.0b'_2 \dots b'_{p-4}110$	GT
$1.0b_2 \dots b_{p-4}011$	$1.1b'_2 \dots b'_{p-4}000$	LT
	$1.1b'_2 \dots b'_{p-4}100$	LT
	$1.0b'_2 \dots b'_{p-4}000$	LT
	$1.0b'_2 \dots b'_{p-4}010$	LT
	$1.0b'_2 \dots b'_{p-4}100$	LT
	$1.0b'_2 \dots b'_{p-4}110$	LT
$1.0b_2 \dots b_{p-4}100$	$1.1b'_2 \dots b'_{p-4}010$	EQ
	$1.1b'_2 \dots b'_{p-4}110$	EQ
	$1.0b'_2 \dots b'_{p-4}001$	EQ
	$1.0b'_2 \dots b'_{p-4}011$	EQ
	$1.0b'_2 \dots b'_{p-4}101$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	EQ
$1.0b_2 \dots b_{p-4}101$	$1.1b'_2 \dots b'_{p-4}000$	GT
	$1.1b'_2 \dots b'_{p-4}100$	GT
	$1.0b'_2 \dots b'_{p-4}000$	GT
	$1.0b'_2 \dots b'_{p-4}010$	GT
	$1.0b'_2 \dots b'_{p-4}100$	GT
	$1.0b'_2 \dots b'_{p-4}110$	GT
$1.0b_2 \dots b_{p-4}110$	$1.1b'_2 \dots b'_{p-4}001$	EQ
	$1.1b'_2 \dots b'_{p-4}101$	EQ
	$1.0b'_2 \dots b'_{p-4}000$	LT
	$1.0b'_2 \dots b'_{p-4}010$	LT
	$1.0b'_2 \dots b'_{p-4}100$	LT
	$1.0b'_2 \dots b'_{p-4}110$	LT
$1.0b_2 \dots b_{p-4}111$	$1.1b'_2 \dots b'_{p-4}010$	LT
	$1.1b'_2 \dots b'_{p-4}110$	LT
	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.0b'_2 \dots b'_{p-4}011$	LT
	$1.0b'_2 \dots b'_{p-4}101$	LT
	$1.0b'_2 \dots b'_{p-4}111$	LT

Table 2: Result and sign of the error for the multiplication by 3 of $x = 1.1b_2 \dots b_{p-1} \times 2^{e_x}$

\mathbf{x}	$\mathbf{fl(3x)}$	Error
$1.1b_2 \dots b_{p-4}000$	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}010$	EQ
	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}110$	EQ
$1.1b_2 \dots b_{p-4}001$	$1.0b'_2 \dots b'_{p-4}001$	GT
	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.0b'_2 \dots b'_{p-4}101$	GT
	$1.0b'_2 \dots b'_{p-4}111$	GT
$1.1b_2 \dots b_{p-4}010$	$1.0b'_2 \dots b'_{p-4}000$	GT
	$1.0b'_2 \dots b'_{p-4}010$	GT
	$1.0b'_2 \dots b'_{p-4}100$	GT
	$1.0b'_2 \dots b'_{p-4}110$	GT
$1.1b_2 \dots b_{p-4}011$	$1.0b'_2 \dots b'_{p-4}000$	LT
	$1.0b'_2 \dots b'_{p-4}010$	LT
	$1.0b'_2 \dots b'_{p-4}100$	LT
	$1.0b'_2 \dots b'_{p-4}110$	LT
$1.1b_2 \dots b_{p-4}100$	$1.0b'_2 \dots b'_{p-4}001$	EQ
	$1.0b'_2 \dots b'_{p-4}011$	EQ
	$1.0b'_2 \dots b'_{p-4}101$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	EQ
$1.1b_2 \dots b_{p-4}101$	$1.0b'_2 \dots b'_{p-4}000$	GT
	$1.0b'_2 \dots b'_{p-4}010$	GT
	$1.0b'_2 \dots b'_{p-4}100$	GT
	$1.0b'_2 \dots b'_{p-4}110$	GT
$1.1b_2 \dots b_{p-4}110$	$1.0b'_2 \dots b'_{p-4}000$	LT
	$1.0b'_2 \dots b'_{p-4}010$	LT
	$1.0b'_2 \dots b'_{p-4}100$	LT
	$1.0b'_2 \dots b'_{p-4}110$	LT
$1.1b_2 \dots b_{p-4}111$	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.0b'_2 \dots b'_{p-4}011$	LT
	$1.0b'_2 \dots b'_{p-4}101$	LT
	$1.0b'_2 \dots b'_{p-4}111$	LT

Proof. Trivial. Consider:

$$\begin{array}{r}
 \times \quad \begin{array}{ccccccc} 0. & 0 & b_2 & \cdots & b_{p-1} & & \times 2^{E_{\min}} \\ 1. & 1 & & & & & \times 2^1 \end{array} \\
 \hline
 + \quad \begin{array}{ccccccc} 0. & 0 & b_2 & b_3 & \cdots & b_{p-1} & \\ 0. & 0 & b_2 & b_3 & \cdots & b_{p-1} & \end{array} \\
 \hline
 \quad \begin{array}{ccccccc} 0 & b'_1 & b'_2 & b'_3 & \cdots & b'_{p-1} & b'_p & \times 2^{E_{\min}} \end{array}
 \end{array} \tag{7}$$

□

On the other hand, if $b_1 \neq 0$, we have:

$$\begin{array}{r}
 \times \quad \begin{array}{ccccccc} 0. & 1 & b_2 & \cdots & b_{p-1} & & \times 2^{E_{\min}} \\ 1. & 1 & & & & & \times 2^1 \end{array} \\
 \hline
 + \quad \begin{array}{ccccccc} 0. & 1 & b_2 & \cdots & b_{p-1} & & \\ 0. & 1 & b_2 & \cdots & b_{p-1} & & \\ c_0 & c_1 & & & & & \end{array} \\
 \hline
 \quad \begin{array}{ccccccc} b'_0 & b'_1 & b'_2 & b'_3 & \cdots & b'_{p-1} & b'_p & \times 2^{E_{\min}+1} \end{array}
 \end{array}$$

where c_0 and c_1 are, as before, the leftmost carries.

It is easy to see that c_0 and c_1 must be equal. Furthermore:

$$\left\{ \begin{array}{l} \text{if } c_0 = c_1 = 0, \quad \text{then } b'_0 = 0 \text{ and } b'_1 = 1 \\ \text{if } c_0 = c_1 = 1, \quad \text{then } b'_0 = 1 \text{ and } b'_1 = 0 \end{array} \right.$$

Multiplying x by 3 can then lead to a p bits significand (when $b'_0 = 0$) or to a $p + 1$ bits significand (when $b'_0 = 1$). Appendix B displays the different results for all possible values of the rightmost bits².

From Appendix B and Eq. (7), we can deduce the value of $\text{fl}(3x)$ obtained from rounding $3x$ to nearest-even (see Table 3).

4 Dividing x by three

Similarly to the multiplication, we will first consider the case in which x is a normal number and no underflow occurs; we will then investigate what happens when x is a denormal number or when the result underflows.

For $x = (-1)^{s_x} m_x \times 2^{e_x}$ and $y = (-1)^{s_y} m_y \times 2^{e_y}$, the quotient x/y is given by [9]:

$$x/y = (-1)^{s_x \oplus s_y} m_x/m_y \times 2^{e_x - e_y} \tag{8}$$

To determine the sign of the error of $\text{fl}(x/y)$ relative to x/y , we may only consider the significands m_x and m_y . Consequently, we will only display the exponents when rendered necessary by the context. We also restrict ourselves to positive values for x , the negative case being easily deduced.

4.1 The normal case

As for the multiplication, we need only consider the first bit after the radix point and the last three bits of the significand to determine the result of the division by 3. The following three lemmas tell us how many bits to compute in order to obtain a p bits significand correctly rounded to nearest-even.

Lemma 4. *If $x = 1.1b_2 \dots b_{p-1} \times 2^{e_x}$ is a normal floating-point number, then we need to compute $p + 1$ bits to get a correctly rounded p bits significand for $y = \text{fl}(x/3)$; in addition, we have:*

$$\text{fl}(y) = 1.0b'_2 \dots b'_{p-1} \times 2^{e_x-1} \tag{9}$$

provided no underflow occurs.

²Note that, even though only the last two bits need to be considered, we take into account the last three as in the previous section.

Table 3: Result and sign of the error for the multiplication by 3 of $x = 0.b_1b_2 \dots b_{p-1} \times 2^{E_{\min}}$

\mathbf{x}	$\mathbf{fl(3x)}$	Error
$0.0b_2 \dots b_{p-1}$	$b'_0.b'_1 \dots b'_{p-1}$	EQ
$0.1b_2 \dots b_{p-4}000$	$1.1b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-3}00$	EQ
$0.1b_2 \dots b_{p-4}001$	$1.1b'_2 \dots b'_{p-4}011$	EQ
	$1.0b'_2 \dots b'_{p-3}10$	GT
$0.1b_2 \dots b_{p-4}010$	$1.1b'_2 \dots b'_{p-4}110$	EQ
	$1.0b'_2 \dots b'_{p-3}110$	EQ
$0.1b_2 \dots b_{p-4}011$	$1.1b'_2 \dots b'_{p-4}001$	EQ
	$1.0b'_2 \dots b'_{p-3}00$	LT
$0.1b_2 \dots b_{p-4}100$	$1.1b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-3}10$	EQ
$0.1b_2 \dots b_{p-4}101$	$1.1b'_2 \dots b'_{p-4}111$	EQ
	$1.0b'_2 \dots b'_{p-3}00$	GT
$0.1b_2 \dots b_{p-4}110$	$1.1b'_2 \dots b'_{p-4}010$	EQ
	$1.0b'_2 \dots b'_{p-3}01$	EQ
$0.1b_2 \dots b_{p-4}111$	$1.1b'_2 \dots b'_{p-4}101$	EQ
	$1.0b'_2 \dots b'_{p-3}10$	LT

Proof. Consider the process of dividing x by 3 at its beginning:

$$\begin{array}{r}
 0 \ 1 \ 0 \ \dots \\
 11 \overline{) 1 \ 1 \ b_2 \ \dots \ b_{p-1}} \\
 \underline{1 \ 1} \\
 0 \ b_2 \\
 \dots
 \end{array}$$

Due to the leading 0 of the quotient, we need to compute one bit more to get p bits in the rounded significand. Equation (9) is obvious from the division process shown above. \square

Lemma 5. *If $x = 1.0b_2 \dots b_{p-1} \times 2^{E_x}$ is a normal floating-point number, then we need to compute $p + 2$ bits to get a p bit significand for $y = x/3$.*

Proof. Consider the process of dividing x by 3 at its beginning:

$$\begin{array}{r}
 0 \ 0 \ 1 \ b'_1 \ \dots \\
 11 \overline{) 1 \ 0 \ b_2 \ b_3 \ \dots \ b_{p-1}} \\
 \underline{1 \ 0} \\
 r_0 \ r_1 \ b_3 \\
 \dots
 \end{array}$$

Due to the two leading 0 of the quotient, we need to compute two bits more to get p bits in the significand of the result. \square

Lemma 6. *Given $x = 1.b_1 \dots b_{p-1} \times 2^{E_x}$ a normal floating-point number and $y = x/3$. There are three possible forms for the significand of y :*

$$\left\{ \begin{array}{l} 1.b'_1 \dots b'_{p-2} \overline{0} \\ 1.b'_1 \dots b'_{p-2} \overline{01} \\ 1.b'_1 \dots b'_{p-2} \overline{10} \end{array} \right.$$

where \bar{s} denotes an infinite repetition of the binary string s .

Proof. In the process of dividing x by 3, there are only three possible remainders at each step: 0_2 , 1_2 , or 10_2 . Consequently, when $x = 1.1b_2 \dots b_{p-1} \times 2^{e_x}$, we only have the following three situations (the case $x = 1.0b_2 \dots b_{p-1} \times 2^{e_x}$ is analogous):

$$\begin{array}{r}
 \begin{array}{cccccccc}
 & 0 & 1 & 0 & b'_2 & \dots & b'_{p-2} & \overline{0} \dots \\
 11) & 1 & 1 & b_2 & \dots & & b_{p-1} & \\
 & & & & & & & \downarrow \\
 & & & & & & \boxed{0} & 0
 \end{array}
 &
 \begin{array}{cccccccc}
 & 0 & 1 & 0 & b'_2 & \dots & b'_{p-2} & \overline{0 \ 1} \dots \\
 11) & 1 & 1 & b_2 & \dots & & b_{p-1} & \\
 & & & & & & & \downarrow \downarrow \\
 & & & & & & \boxed{1} & 0 \ 0 \\
 & & & & & & & 1
 \end{array} \\
 \\
 \begin{array}{cccccccc}
 & 0 & 1 & 0 & b'_2 & \dots & b'_{p-2} & \overline{1 \ 0} \dots \\
 11) & 1 & 1 & b_2 & \dots & & b_{p-1} & \\
 & & & & & & & \downarrow \\
 & & & & & & \boxed{10} & 0 \ \downarrow \\
 & & & & & & & 1 \ 0
 \end{array}
 \end{array} \tag{10}$$

□

For each triplet of rightmost bits of the significand of x , there are three cases to consider depending on the remainder after having taken into account the bit b_{p-4} . All the cases are presented in Section C.

Tables 4 and 5 summarize the results for the whole section.

4.2 The denormal case

There are two possibilities for an underflow to occur when dividing x by 3: either x is already a denormal number, or x is normal but the rounded result of the division is a denormal number.

If x is a denormal number, we have:

$$\begin{array}{r}
 \begin{array}{cccccccc}
 & 0. & 0 & b'_2 & \dots & & b'_{p-1} \times 2^{E_{\min}} \\
 1.1 \times 2^1) & 0. & b_1 & \dots & b_{p-4} & b_{p-3} & b_{p-2} & b_{p-1} \times 2^{E_{\min}} \\
 & & & \dots & & & & \\
 & & & & \dots & & & \\
 & & & & & \dots & & \\
 & & & & & & \dots & \\
 & & & & & & & \dots
 \end{array}
 \end{array} \tag{11}$$

If x is a normal number of the form $x = 1.b_1 \dots b_{p-1} \times 2^{E_{\min}}$, we have, either:

$$\begin{array}{r}
 \begin{array}{cccccccc}
 & 0. & 0 & b'_2 & \dots & & b'_{p-1} \times 2^{E_{\min}} \\
 1.1 \times 2^1) & 1. & 0 & \dots & b_{p-4} & b_{p-3} & b_{p-2} & b_{p-1} \times 2^{E_{\min}} \\
 & & & \dots & & & & \\
 & & & & \dots & & & \\
 & & & & & \dots & & \\
 & & & & & & \dots & \\
 & & & & & & & \dots
 \end{array}
 \end{array} \tag{12}$$

or:

$$\begin{array}{r}
 \begin{array}{cccccccc}
 & 0. & 1 & b'_2 & \dots & & b'_{p-1} \times 2^{E_{\min}} \\
 1.1 \times 2^1) & 1. & 1 & \dots & b_{p-4} & b_{p-3} & b_{p-2} & b_{p-1} \times 2^{E_{\min}} \\
 & & & \dots & & & & \\
 & & & & \dots & & & \\
 & & & & & \dots & & \\
 & & & & & & \dots & \\
 & & & & & & & \dots
 \end{array}
 \end{array} \tag{13}$$

In all three cases, the result of the division cannot be scaled to the right since the exponent is the smallest possible. As a consequence, we need only compute p bits to obtain $\text{fl}(x/3) = 0.b'_1 \dots b'_{p-3} b'_{p-2} b'_{p-1} \times 2^{E_{\min}}$.

Table 4: Result and sign of the error for the division by 3 of $x = 1.0b_2 \dots b_{p-1} \times 2^{e_x}$ with no underflow

\mathbf{x}	$\text{fl}(\mathbf{x}/3)$	Error
$1.0b_2 \dots b_{p-4}000$	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.1b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.1b'_2 \dots b'_{p-4}011$	GT
	$1.0b'_2 \dots b'_{p-4}101$	LT
	$1.1b'_2 \dots b'_{p-4}101$	LT
$1.0b_2 \dots b_{p-4}001$	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.1b'_2 \dots b'_{p-4}001$	LT
	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.1b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	GT
	$1.1b'_2 \dots b'_{p-4}111$	GT
$1.0b_2 \dots b_{p-4}010$	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.1b'_2 \dots b'_{p-4}011$	GT
	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.1b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}101$	LT
	$1.1b'_2 \dots b'_{p-4}101$	LT
$1.0b_2 \dots b_{p-4}011$	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.1b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	GT
	$1.1b'_2 \dots b'_{p-4}111$	GT
	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.1b'_2 \dots b'_{p-4}001$	LT
$1.0b_2 \dots b_{p-4}100$	$1.0b'_2 \dots b'_{p-4}101$	LT
	$1.1b'_2 \dots b'_{p-4}101$	LT
	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.1b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.1b'_2 \dots b'_{p-4}011$	GT
$1.0b_2 \dots b_{p-4}101$	$1.0b'_2 \dots b'_{p-4}111$	GT
	$1.1b'_2 \dots b'_{p-4}111$	GT
	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.1b'_2 \dots b'_{p-4}001$	LT
	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.1b'_2 \dots b'_{p-4}100$	EQ
$1.0b_2 \dots b_{p-4}110$	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.1b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.1b'_2 \dots b'_{p-4}011$	GT
	$1.0b'_2 \dots b'_{p-4}101$	LT
	$1.1b'_2 \dots b'_{p-4}101$	LT
$1.0b_2 \dots b_{p-4}111$	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.1b'_2 \dots b'_{p-4}001$	LT
	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.1b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	GT
	$1.1b'_2 \dots b'_{p-4}111$	GT

Table 5: Result and sign of the error for the division by 3 of $x = 1.1b_2 \dots b_{p-1} \times 2^{e_x}$ with no underflow.

\mathbf{x}	$\text{fl}(\mathbf{x}/3)$	Error
$1.1b_2 \dots b_{p-4}000$	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}101$	LT
	$1.0b'_2 \dots b'_{p-4}011$	GT
$1.1b_2 \dots b_{p-4}001$	$1.0b'_2 \dots b'_{p-4}001$	GT
	$1.0b'_2 \dots b'_{p-4}110$	EQ
	$1.0b'_2 \dots b'_{p-4}011$	LT
$1.1b_2 \dots b_{p-4}010$	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	GT
$1.1b_2 \dots b_{p-4}011$	$1.0b'_2 \dots b'_{p-4}010$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	LT
	$1.0b'_2 \dots b'_{p-4}101$	GT
$1.1b_2 \dots b_{p-4}100$	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.0b'_2 \dots b'_{p-4}000$	EQ
	$1.0b'_2 \dots b'_{p-4}101$	LT
$1.1b_2 \dots b_{p-4}101$	$1.0b'_2 \dots b'_{p-4}011$	LT
	$1.0b'_2 \dots b'_{p-4}001$	GT
	$1.0b'_2 \dots b'_{p-4}110$	EQ
$1.1b_2 \dots b_{p-4}110$	$1.0b'_2 \dots b'_{p-4}100$	EQ
	$1.0b'_2 \dots b'_{p-4}001$	LT
	$1.0b'_2 \dots b'_{p-4}111$	GT
$1.1b_2 \dots b_{p-4}111$	$1.0b'_2 \dots b'_{p-4}101$	GT
	$1.0b'_2 \dots b'_{p-4}010$	EQ
	$1.0b'_2 \dots b'_{p-4}111$	LT

Table 7: Result and sign of the error for the division by 3 of $x = 1.0b_2 \dots b_{p-1} \times 2^{E_{\min}+1}$ with underflow

\mathbf{x}	$\text{fl}(\mathbf{x}/3)$	Error
$1.0b_2 \dots b_{p-4}000$	$0.1 \dots b'_{p-4}000$	EQ
	$0.1 \dots b'_{p-4}101$	LT
	$0.1 \dots b'_{p-4}011$	GT
$1.0b_2 \dots b_{p-4}001$	$0.1 \dots b'_{p-4}001$	GT
	$0.1 \dots b'_{p-4}110$	EQ
	$0.1 \dots b'_{p-4}011$	LT
$1.0b_2 \dots b_{p-4}010$	$0.1 \dots b'_{p-4}001$	LT
	$0.1 \dots b'_{p-4}100$	EQ
	$0.1 \dots b'_{p-4}111$	GT
$1.0b_2 \dots b_{p-4}011$	$0.1 \dots b'_{p-4}010$	EQ
	$0.1 \dots b'_{p-4}111$	LT
	$0.1 \dots b'_{p-4}101$	GT
$1.0b_2 \dots b_{p-4}100$	$0.1 \dots b'_{p-4}011$	GT
	$0.1 \dots b'_{p-4}000$	EQ
	$0.1 \dots b'_{p-4}101$	LT
$1.0b_2 \dots b_{p-4}101$	$0.1 \dots b'_{p-4}011$	LT
	$0.1 \dots b'_{p-4}001$	GT
	$0.1 \dots b'_{p-4}110$	EQ
$1.0b_2 \dots b_{p-4}110$	$0.1 \dots b'_{p-4}100$	EQ
	$0.1 \dots b'_{p-4}001$	LT
	$0.1 \dots b'_{p-4}111$	GT
$1.0b_2 \dots b_{p-4}111$	$0.1 \dots b'_{p-4}101$	GT
	$0.1 \dots b'_{p-4}010$	EQ
	$0.1 \dots b'_{p-4}111$	LT

5 Composing the multiplication and division

As seen in Section 2, the correctness of the floating-point multiplication and division are precisely defined by IEEE 754. From Eq. (2) and Eq. (3), we may compute the worst case error for the composition of a multiplication and a division. In theory, there are eight cases depending on the order of the operations and on the outcome of each operation (underflow or not). Tables 8 and 9 summarize the error analysis when computing either $\text{fl}\langle(3x)/3\rangle$ or $\text{fl}\langle 3(x/3)\rangle$. The crossed out cells in the tables correspond to cases that cannot occur: Consider, for example, Table 8; if there is an underflow when multiplying x by 3, there will surely be an underflow when dividing $3x$ by 3.

Table 8: Worst case error analysis for $\text{fl}\langle\frac{3x}{3}\rangle$, with: $|\eta_1| \leq \mu/2$, $|\eta_2| \leq \mu/2$, $|\delta_1| \leq u/(1+u)$, $|\delta_2| \leq u - 2u^2$.

		$z = y/3$	
		Underflow	No underflow
$y = 3x$	Underflow	$\frac{3x + \eta_1}{3} + \eta_2$	$\frac{3x + \eta_1}{3}(1 + \delta_2)$
	No underflow	$\frac{3x(1 + \delta_1)}{3} + \eta_2$	$\frac{3x}{3}(1 + \delta_1)(1 + \delta_2)$

Table 9: Worst case error analysis for $\text{fl}\langle 3\frac{x}{3}\rangle$, with: $|\eta_1| \leq \mu/2$, $|\eta_2| \leq \mu/2$, $|\delta_1| \leq u - 2u^2$, $|\delta_2| \leq u/(1+u)$.

		$z = 3y$	
		Underflow	No underflow
$y = x/3$	Underflow	$3\left(\frac{x}{3} + \eta_1\right) + \eta_2$	$3\left(\frac{x}{3} + \eta_1\right)(1 + \delta_2)$
	No underflow	$3\left(\frac{x}{3}(1 + \delta_1)\right) + \eta_2$	$3\frac{x}{3}(1 + \delta_1)(1 + \delta_2)$

5.1 Multiplication and division with no underflow

Considering first the case in which no underflow occurs, we get, irrespective of the order of the operations considered:

$$\text{fl}\langle\frac{3x}{3}\rangle = \text{fl}\langle 3\frac{x}{3}\rangle = x(1 + \delta_a)(1 + \delta_b)$$

with $|\delta_a| \leq u/(1+u)$ and $|\delta_b| \leq u - 2u^2$. Equivalently, we have:

$$\text{fl}\langle\frac{3x}{3}\rangle = \text{fl}\langle 3\frac{x}{3}\rangle = x(1 + \delta) \quad (15)$$

with $\delta = \delta_a + \delta_b + \delta_a\delta_b$, which means that $|\delta| \leq 2(u - 2u^3)/(1+u)$, that is $|\delta| < 2u$. The relative distance between two consecutive floating-point numbers is contained in the range $[u, 2u]$; Consequently, the rounded value of $3x/3$ or $3(x/3)$ may not round to x exactly, but we have the guarantee that it must be at worst one of the immediate floating-point neighbors of x (See Fig. 1).

5.2 Multiplication and division with underflow

As can be seen in Table 8 and Table 9, the worst case errors are different for $\text{fl}\langle(3x)/3\rangle$ and $\text{fl}\langle 3(x/3)\rangle$ whenever an underflow occurs. We must, therefore, consider separately the two computations.

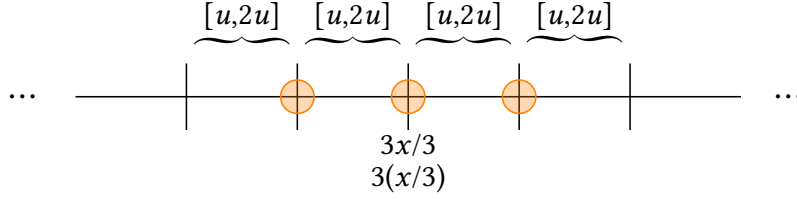


Figure 1: Orange disks representing the possible values for $\text{fl}\langle 3x/3 \rangle$ or $\text{fl}\langle 3(x/3) \rangle$ with rounding to nearest-even in the absence of underflow.

5.2.1 The case of $\text{fl}\langle (3x)/3 \rangle$

Underflow on multiplication and division. When both the multiplication and division underflow, we get:

$$\text{fl}\langle \frac{3x}{3} \rangle = \frac{3x + \eta_1}{3} + \eta_2 = x + \left(\frac{\eta_1}{3} + \eta_2 \right), \text{ with } |\eta_1| \leq \frac{\mu}{2}, |\eta_2| \leq \frac{\mu}{2}$$

Hence,

$$\text{fl}\langle \frac{3x}{3} \rangle = x + \eta, \text{ with } |\eta| \leq \frac{2\mu}{3} \quad (16)$$

For the multiplication to underflow, x must be a subnormal number. As shown in Figure 2, the absolute distance between two consecutive subnormals is μ . Consequently, we deduce from Eq. (16) that

$$\text{fl}\langle \frac{3x}{3} \rangle = \frac{3x}{3} = x$$

if both the multiplication and the division underflow.

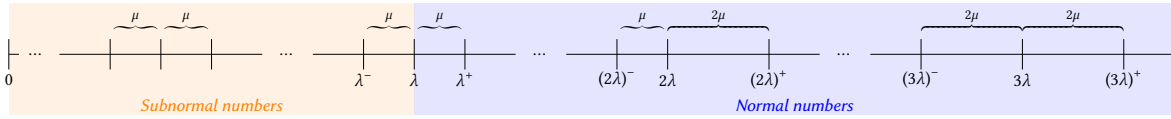


Figure 2: Absolute distance between small positive floating-point numbers

Underflow on division but not on multiplication. If the multiplication does not underflow but the division does, we get:

$$\text{fl}\langle \frac{3x}{3} \rangle = \frac{3x(1 + \delta_1)}{3} + \eta_2, \text{ with } |\delta_1| \leq \frac{u}{1 + u}, |\eta_2| \leq \frac{\mu}{2}.$$

It is easy to check that for $x = \lambda$, we have $\text{fl}\langle (3x)/3 \rangle = \lambda$. By monotonicity of rounding, we get that for all y greater than x , we have $\text{fl}\langle (3y)/3 \rangle \geq \lambda$. Hence, x must be a subnormal number for the division to underflow when computing $(3x)/3$. Therefore, the absolute distance between $3x$ and $\text{fl}\langle 3x \rangle$ must be smaller or equal to μ (See Figure 2), which leads to:

$$\begin{aligned} \text{fl}\langle \frac{3x}{3} \rangle &= \frac{3x + \eta_3}{3} + \eta_2, \text{ with } |\eta_3| \leq \mu, |\eta_2| \leq \frac{\mu}{2} \\ &= x + \frac{5}{6}\mu \end{aligned}$$

As in the previous case, we deduce:

$$\text{fl}\langle \frac{3x}{3} \rangle = \frac{3x}{3} = x$$

if the multiplication does not underflow but the division does.

5.2.2 The case of $\text{fl}\langle 3(x/3) \rangle$

Underflow on division and multiplication. If both the division and the multiplication underflow, we get:

$$\text{fl}\langle 3\frac{x}{3} \rangle = 3 \left(\frac{x}{3} + \eta_1 \right) + \eta_2, \text{ with } |\eta_1| \leq \mu/2, |\eta_2| \leq \mu/2 \quad (17)$$

Using monotonicity of rounding, if $\text{fl}\langle 3(x/3) \rangle$ underflows, we must have $x \leq \lambda$. Therefore, $\text{fl}(x/3)$ is of the form $0.0b'_2 \dots b'_{p-1} \times 2^{E_{\min}}$. According to Table 3, the multiplication is then without error, that is: $\eta_2 = 0$. Equation (17) simplifies to:

$$\text{fl}\langle 3\frac{x}{3} \rangle = x + 3\eta_1, \text{ with } |\eta_1| \leq \mu/2$$

Hence,

$$|\text{fl}\langle 3\frac{x}{3} \rangle - x| \leq \frac{3\mu}{2}$$

Since the absolute distance between subnormal numbers is μ , we have:

$$\text{fl}\langle 3\frac{x}{3} \rangle \in \{x^-, x, x^+\}$$

Underflow on division but not on multiplication. If the division underflows but the multiplication does not, we get:

$$\text{fl}\langle 3\frac{x}{3} \rangle = 3 \left(\frac{x}{3} + \eta_1 \right) (1 + \delta_2), \text{ with } |\eta_1| \leq \mu/2, |\delta_2| \leq u/(1+u) \quad (18)$$

It is trivial to check that we must have $x \in [\lambda^+, (3\lambda)^-]$, using monotonicity of rounding. For convenience, we will investigate the actual worst case error on three disjoint domains for x :

- If $x \in [\lambda^+, (\frac{3}{2}\lambda)^-]$ (i.e., $x \in [1.00 \dots 001 \times 2^{E_{\min}}, 1.011 \dots 111 \times 2^{E_{\min}}]$).

According to Table 6, we have:

$$x \in [\lambda^+, (\frac{3}{2}\lambda)^-] \implies \text{fl}\left(\frac{x}{3}\right) = 0.0b'_2 \dots b'_{p-1} \times 2^{E_{\min}}$$

According to Table 3, we therefore have:

$$\text{fl}\left(3\text{fl}\left(\frac{x}{3}\right)\right) = 3\text{fl}\left(\frac{x}{3}\right)$$

which means:

$$\text{fl}\langle 3\frac{x}{3} \rangle = 3\left(\frac{x}{3} + \eta_1\right), |\eta_1| \leq \mu/2$$

Therefore:

$$|\text{fl}\langle 3\frac{x}{3} \rangle - x| \leq \frac{3\mu}{2}$$

Since, the absolute distance between floating-point numbers is μ in the domain $[\lambda^+, (\frac{3}{2}\lambda)^-]$, we have:

$$\forall x \in [\lambda^+, (3\lambda/2)^-] : \text{fl}\langle 3\frac{x}{3} \rangle \in \{x^-, x, x^+\} \quad (19)$$

- If $x \in [3\lambda/2, (2\lambda)^-]$ (i.e. $x \in [1.10 \dots 000 \times 2^{E_{\min}}, 1.11 \dots 111 \times 2^{E_{\min}}]$)

It is easy to check, using monotonicity of rounding, that:

$$x \in [3\lambda/2, (2\lambda)^-] \implies \text{fl}\langle 3\frac{x}{3} \rangle \in [3\lambda/2, ((2\lambda)^-)^-]$$

From Table 3, we get for y a subnormal number:

$$\text{fl}(3y) = 1.1b'_2 \dots b'_{p-1} \times 2^{E_{\min}} \implies \text{fl}(3y) = 3y$$

Hence:

$$\forall x \in [3\lambda/2, (2\lambda)^-]: \text{fl}\langle 3\frac{x}{3} \rangle = 3(\frac{x}{3} + \eta_1) = x + 3\eta_1, |\eta_1| \leq \mu/2$$

From which we get:

$$|\text{fl}\langle 3\frac{x}{3} \rangle - x| \leq \frac{3\mu}{2}$$

Consequently, we have:

$$\forall x \in [3\lambda/2, (2\lambda)^-]: \text{fl}\langle 3\frac{x}{3} \rangle \in \{x^-, x, x^+\} \quad (20)$$

- If $x \in [2\lambda, (3\lambda)^-]$ (i.e., $x \in [1.0 \times 2^{E_{\min}+1}, 1.01 \dots 111 \times 2^{E_{\min}+1}]$).

It is trivial to show, using monotonicity of rounding, that:

$$x \in [2\lambda, (3\lambda)^-] \implies \text{fl}\langle 3\frac{x}{3} \rangle \in [2\lambda, ((3\lambda)^-)^-]$$

In the interval $[2\lambda, ((3\lambda)^-)^-]$, the absolute distance between two consecutive floating-point numbers is 2μ (See Figure 2), which means that:

$$\forall y: \text{fl}(3y) \in [2\lambda, ((3\lambda)^-)^-] \implies |\text{fl}(3y) - 3y| \leq \eta_3, |\eta_3| \leq \mu$$

Therefore, for $x \in [2\lambda, (3\lambda)^-]$:

$$\begin{aligned} \text{fl}\langle 3\frac{x}{3} \rangle &= 3(\frac{x}{3} + \eta_1) + \eta_3, |\eta_1| \leq \mu/2, |\eta_3| \leq \mu \\ &= x + \eta, |\eta| \leq 5\mu/2 \end{aligned}$$

Then:

$$|\text{fl}\langle 3\frac{x}{3} \rangle - x| \leq \frac{5\mu}{2}$$

The distance between consecutive floating-point numbers in the domain involved being 2μ , we eventually get:

$$\forall x \in [2\lambda, (3\lambda)^-]: \text{fl}\langle 3\frac{x}{3} \rangle \in \{x^-, x, x^+\} \quad (21)$$

From Equations (19), (20), (21), we deduce that for $x \in [\lambda^+, (3\lambda)^-]$, $\text{fl}\langle 3(x/3) \rangle$ is equal to x itself or one of its immediate predecessor or successor:

$$\forall x \in [\lambda^+, (3\lambda)^-]: \text{fl}\langle 3\frac{x}{3} \rangle \in \{x^-, x, x^+\} \quad (22)$$

5.3 Putting it all together

Knowing the worst case error for each operation, and armed with the tables from Sections 3 and 4, we may now determine the possible signs of the error when computing $\text{fl}\langle (3x)/3 \rangle$ and $\text{fl}\langle 3(x/3) \rangle$. Consider for example the floating-point number x of the form $1.1b_2 \dots b_{p-4}001 \times 2^{E_x}$ and the expression $(3x)/3$. Using Table 2, we get all the possible forms for the significand of $y = \text{fl}(3x)$. They are—in the absence of underflow:

$$\left\{ \begin{array}{l} 1.0b'_2 \dots b'_{p-4}001 \\ 1.0b'_2 \dots b'_{p-4}011 \\ 1.0b'_2 \dots b'_{p-4}101 \\ 1.0b'_2 \dots b'_{p-4}111 \end{array} \right. \quad (23)$$

Table 10: Result and sign of the error when computing $\text{fl}\langle(3x)/3\rangle$ without underflow.

\mathbf{x}	$\text{fl}\langle(3\mathbf{x})/3\rangle$	Error
$1.0b_2 \dots b_{p-4}000$	$1.0b'_2 \dots b'_{p-4}000$	EQ
$1.0b_2 \dots b_{p-4}001$	$1.0b'_2 \dots b'_{p-4}001$	EQ
$1.0b_2 \dots b_{p-4}010$	$1.0b'_2 \dots b'_{p-4}011$	GT
	$1.0b'_2 \dots b'_{p-4}010$	EQ
$1.0b_2 \dots b_{p-4}011$	$1.0b'_2 \dots b'_{p-4}011$	EQ
$1.0b_2 \dots b_{p-4}100$	$1.0b'_2 \dots b'_{p-4}100$	EQ
$1.0b_2 \dots b_{p-4}101$	$1.0b'_2 \dots b'_{p-4}101$	EQ
$1.0b_2 \dots b_{p-4}110$	$1.0b'_2 \dots b'_{p-4}110$	EQ
	$1.0b'_2 \dots b'_{p-4}101$	LT
$1.0b_2 \dots b_{p-4}111$	$1.0b'_2 \dots b'_{p-4}111$	EQ
$1.1b_2 \dots b_{p-4}000$	$1.1b'_2 \dots b'_{p-4}000$	EQ
$1.1b_2 \dots b_{p-4}001$	$1.1b'_2 \dots b'_{p-4}001$	EQ
$1.1b_2 \dots b_{p-4}010$	$1.1b'_2 \dots b'_{p-4}011$	GT
$1.1b_2 \dots b_{p-4}011$	$1.1b'_2 \dots b'_{p-4}011$	EQ
$1.1b_2 \dots b_{p-4}100$	$1.1b'_2 \dots b'_{p-4}100$	EQ
$1.1b_2 \dots b_{p-4}101$	$1.1b'_2 \dots b'_{p-4}101$	EQ
$1.1b_2 \dots b_{p-4}110$	$1.1b'_2 \dots b'_{p-4}101$	LT
$1.1b_2 \dots b_{p-4}111$	$1.1b'_2 \dots b'_{p-4}111$	EQ

For each of these four forms, we get from Table 4 all the possible forms for the significand of $\text{fl}\langle(y/3)\rangle$. For example, for the first form $1.0b'_2 \dots b'_{p-4}001$, we get:

$$\left\{ \begin{array}{l} 1.0b''_2 \dots b''_{p-4}001 \\ 1.1b''_2 \dots b''_{p-4}001 \\ 1.0b''_2 \dots b''_{p-4}100 \\ 1.1b''_2 \dots b''_{p-4}100 \\ 1.0b''_2 \dots b''_{p-4}111 \\ 1.1b''_2 \dots b''_{p-4}111 \end{array} \right. \quad (24)$$

However, from Eq. (15), we know that $\text{fl}\langle(3x)/3\rangle$ for $x = 1.1b_2 \dots b_{p-4}001 \times 2^{e_x}$ must be one of the following floating-point numbers:

$$\left\{ \begin{array}{l} 1.1b_2 \dots b_{p-4}000 \times 2^{e_x} \\ 1.1b_2 \dots b_{p-4}001 \times 2^{e_x} \\ 1.1b_2 \dots b_{p-4}010 \times 2^{e_x} \end{array} \right.$$

We can then discard from Eq. (24) all forms but the second one: $1.1b_2 \dots b_{p-4}001$. If we do the same work for all forms in Eq. (23), we will find that the only possible significand for $\text{fl}\langle(3x)/3\rangle$ is $1.1b_2 \dots b_{p-4}001$, that is the one of x itself.

Doing the work exemplified above for all possible significands for x that do not lead to an underflow in either the division or the multiplication, we get Table 10. As proved in Section 5.2.1, should an underflow arise during the computation of $\text{fl}\langle(3x)/3\rangle$, it is not necessary to compute such a table since we then always have $\text{fl}\langle(3x)/3\rangle = x$.

Computing the same table for the expression $\text{fl}\langle 3(x/3)\rangle$ results in a table with 96—that is, $2^5 \times 3$ —entries: given $x = b_0.b_1b_2 \dots b_{p-3}b_{p-2}b_{p-1} \times 2^{e_x}$, $\text{fl}\langle 3(x/3)\rangle$ may always be x or one of its two immediate neighbours, whatever the values of b_0 , b_1 , b_{p-3} , b_{p-2} , and b_{p-1} . Therefore, it is useless to display the resulting table as it is uninformative.

6 Conclusion

Using both traditional error analysis and the tables computed in Sections 3 and 4, we have been able to compute refined error bounds in Section 5 for $\text{fl}\langle 3x/3 \rangle$ and $\text{fl}\langle 3(x/3) \rangle$. At first sight (See Tables 11 and 12), it seems that the expression $3x/3$ gives overall more accurate results than the expression $3(x/3)$ whenever an underflow occurs, and that no expression is better than the other when no underflow occurs.

Table 11: Refined worst case error analysis for $\text{fl}\langle \frac{3x}{3} \rangle$, with $|\delta| < 2u$.

$y = 3x$ \ $z = y/3$	Underflow	No underflow
Underflow	x	—
No underflow	x	$x(1 + \delta)$

Table 12: Refined worst case error analysis for $\text{fl}\langle 3\frac{x}{3} \rangle$, with: $|\eta| \leq \mu$, $|\delta| < 2u$.

$y = x/3$ \ $z = 3y$	Underflow	No underflow
Underflow	$\{x^-, x, x^+\}$	$\{x^-, x, x^+\}$
No underflow	—	$x(1 + \delta)$

This is only part of the story, however. Our detailed analysis reveals that only when the last three bits of the significand of x are “010” or “110” may the expression $3x/3$ be rounded to a value different from x (see Table 10), while the same situation may occur whatever the value of x ’s significand with the expression $3(x/3)$ (refer to the end of the previous section). The expression $3x/3$ is therefore to be preferred to $3(x/3)$ in all cases.

Tables 1, 2, 3, 4, 5, 6, and 7 are all interesting and useful in their own right whenever one needs to obtain more information on the division and multiplication by three of a floating-point number than merely a worst-case error bound. We expect to reuse the method used here to prove properties regarding the expressions $3x/3$ and $3(x/3)$, which combines traditional error analysis with systematic study of the behaviour of equivalence classes of floating-point numbers, in proving the correctness of polyadic splitting operators for interval analysis [5].

The question of multiplying and dividing a number by another in floating-point arithmetic has been addressed several times already in slightly different settings: Kahan [3, thm. 7] showed that $\text{fl}\langle (m/n) \times n \rangle = m$ for integers m and n such that $|m| < 2^{p-1}$ and n is the sum of two powers of 2. Edelman [2] investigated the problem of finding the smallest positive floating-point x such that $x(1/x) \neq 1$.

Obviously, the next step for us is to address the problem presented in this paper for other integer values than 3. That might prove indeed necessary for our work on polyadic splitting operators. Once the preliminary worst case error analysis of an expression is performed manually, a large part of the proofs can be performed programmatically, which departs from the more involved—though more powerful and flexible—approaches used by, e.g., Kahan and Edelman. Generalizing the approach to other similar problems of interest is a direction of future research.

References

- [1] Florent De Dinechin and Laurent-Stéphane Didier. “Table-Based Division by Small Integer Constants”. In: *Reconfigurable Computing: Architectures, Tools and Applications*. Ed. by Oliver C. S. Choy et al. Vol. 7199. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 2012, pp. 53–63. ISBN: 978-3-642-28364-2. DOI: [10.1007/978-3-642-28365-9_5](https://doi.org/10.1007/978-3-642-28365-9_5).
- [2] Alan Edelman. “When is $x \times (1/x) \neq 1$?” Unpublished note. Dec. 1994.
- [3] David Goldberg. “What every computer scientist should know about floating-point arithmetic”. In: *ACM Computing Surveys* 23.1 (Mar. 1991), pp. 5–48. ISSN: 0360-0300. DOI: [10.1145/103162.103163](https://doi.org/10.1145/103162.103163).
- [4] Frédéric Goualard. “How do you compute the midpoint of an interval?” In: *ACM Transactions on Mathematical Software* 40.2 (Feb. 2014). DOI: [10.1145/2493882](https://doi.org/10.1145/2493882).
- [5] Frédéric Goualard and Laurent Granvilliers. “Polyadic Splitting of Floating-point Intervals”. In preparation. 2017.
- [6] Nicholas J. Higham. *Accuracy and Stability of Numerical Algorithms*. Second edition. Society for Industrial and Applied Mathematics, 2002. ISBN: 0898715210.
- [7] *IEEE Standard for Binary Floating-Point Arithmetic*. American National Standard (ANSI) IEEE Std 754-1985. The Institute of Electrical and Electronics Engineers, Inc, 1985.
- [8] Claude-Pierre Jeannerod and Siegfried M. Rump. “On relative errors of floating-point operations: optimal bounds and applications”. In: *Mathematics of Computation* (2017). To appear.
- [9] Jean-Michel Muller et al. *Handbook of Floating-Point Arithmetic*. Birkhauser Boston Inc, Dec. 2008. ISBN: 978-0817647049.

A Multiplication by 3 involving no denormal number

The two columns below present the results of the multiplication by 3 of a positive floating-point number x depending on the last three bits of its significand. Thanks to Lemmas 1 and 2, we know that we also need to consider the first bit after the radix point to know how many bits from the result to discard. The left column (resp. right column) corresponds to the cases for which we need to discard the last bit (resp. the last two bits) underlined. For convenience, we display only the significands and leave the exponents implicit.

$$\begin{array}{r}
 \times \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 0 \\
 \hline
 1. \ 1 \ b'_2 \ \dots \ b'_{p-4} \ b'_{p-3} \ 0 \ 0 \ \underline{0}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 1 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 1 \\
 \hline
 1. \ 1 \ b'_2 \ \dots \ b'_{p-4} \ b'_{p-3} \ 0 \ 1 \ \underline{1}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 0 \\
 \hline
 1. \ 1 \ b'_2 \ \dots \ b'_{p-4} \ b'_{p-3} \ 1 \ 1 \ \underline{0}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 1 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 1 \\
 \hline
 1. \ 1 \ b'_2 \ \dots \ b'_{p-4} \ b'_{p-3} \ 0 \ 0 \ \underline{1}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 0 \\
 \hline
 1. \ 1 \ b'_2 \ \dots \ b'_{p-4} \ b'_{p-3} \ 1 \ 0 \ \underline{0}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 1 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 1 \\
 \hline
 1. \ 1 \ b'_2 \ \dots \ b'_{p-4} \ b'_{p-3} \ 1 \ 1 \ \underline{1}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 1 \ 1 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ 0 \ b_2 \ \dots \ b_{p-4} \ 1 \ 1 \ 0 \\
 \hline
 1. \ 1 \ b'_2 \ \dots \ b'_{p-4} \ b'_{p-3} \ 0 \ 1 \ \underline{0}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 0 \\
 \hline
 1. \ 0 \ b'_2 \ \dots \ b'_{p-3} \ b'_{p-2} \ 0 \ \underline{0} \ \underline{0}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 1 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 0 \ 1 \\
 \hline
 1. \ 0 \ b'_2 \ \dots \ b'_{p-3} \ b'_{p-2} \ 0 \ \underline{1} \ \underline{1}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 0 \\
 \hline
 1. \ 0 \ b'_2 \ \dots \ b'_{p-3} \ b'_{p-2} \ 1 \ \underline{1} \ \underline{0}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 1 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 0 \ 1 \ 1 \\
 \hline
 1. \ 0 \ b'_2 \ \dots \ b'_{p-3} \ b'_{p-2} \ 0 \ \underline{0} \ \underline{1}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 0 \\
 \hline
 1. \ 0 \ b'_2 \ \dots \ b'_{p-3} \ b'_{p-2} \ 1 \ \underline{0} \ \underline{0}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 1 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 1 \ 0 \ 1 \\
 \hline
 1. \ 0 \ b'_2 \ \dots \ b'_{p-3} \ b'_{p-2} \ 1 \ \underline{1} \ \underline{1}
 \end{array}$$

$$\begin{array}{r}
 \times \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 1 \ 1 \ 0 \\
 \hline
 \times \quad 1. \ 1 \\
 \hline
 + \quad 1. \ b_1 \ b_2 \ \dots \ b_{p-4} \ 1 \ 1 \ 0 \\
 \hline
 1. \ 0 \ b'_2 \ \dots \ b'_{p-3} \ b'_{p-2} \ 0 \ \underline{1} \ \underline{0}
 \end{array}$$

D Division by 3 involving denormal numbers

As explained in Section 4.2, the rules for the division involving denormal numbers may be inferred from the ones for the division involving normal numbers (Refer to Section C above).

Rapport de recherche Research Report

**When $(3x)/3$ and $3(x/3)$
are not equal to x**

RR n° xx.xx

Frédéric GOUALARD