

# Exact Real Computer Arithmetic

Peter John Potts

Abbas Edalat

Departmental Technical Report DOC 97/9

Imperial College

180 Queen's Gate, London SW7 2BZ, United Kingdom

{pjp,ae}@doc.ic.ac.uk

21 March 1997

## 1 Introduction

Real numbers are usually represented by finite strings of digits belonging to some digit set. The real number representation specifies a function that maps strings to real numbers or real intervals with distinct end-points. For example, IEEE 754 single precision floating point [9] is encoded in 32 binary bits using 1 bit for the sign  $s$ , 8 bits for the biased exponent  $e$ , and 23 bits for the normalised mantissa  $m$  without the leading 1. The basic format represents the real number

$$(-1)^s 2^{e-127} (1.m).$$

However, finite strings of digits can only represent a limited subset of the real numbers exactly because many real numbers have too many significant digits (such as  $\pi$  or  $\sqrt{2}$ ) or are too large or too small. This means that most real numbers are represented by nearby real numbers or enclosing real intervals with distinct end-points giving rise to the notion of round-off errors. This is generally accepted for a wide range of applications. However, it is well-known that the accumulation of round-off errors due to a large number of calculations can produce grossly inaccurate or even incorrect results. Furthermore, current floating point representations in computers have so many bits that verification of floating point functions cannot be achieved by exhaustively testing every possible input combination.

Interval analysis [17] has been used to partially circumvent this problem by maintaining a pair of bounding numbers that is guaranteed to contain the real number or interval in question. However, this interval can get unjustifiably large and thereby convey very little information.

Alternatively, by allowing infinite strings of digits all the real numbers can be represented exactly. The digits in an infinite string are normally used to construct a sequence

of nested real intervals whose lengths converge to zero. The intersection of these intervals is a singleton set whose element is the real number being represented.

Furthermore, basic arithmetic operations are only computable if the representation is *redundant*. In other words, there must be more than one representation for every real number.

In the literature, there are broadly speaking three frameworks for exact real computer arithmetic:

- (i) *Infinite sequences of linear maps* proposed by Avizienis [1] and appeared in the work of Watanuki et al [23], Boehm and Cartwright [2], Di Gianantonio [5], Escardo [4], Nielsen et al [18] and Menissier-Morain [16].
- (ii) *Continued fraction expansions* proposed by Gosper [7], developed by Peyton Jones [10] and Vuillemin [21] and advanced more recently by Kornerup et al [15, 13, 12, 14].
- (iii) *Infinite composition of linear fractional transformations* (also known as homographies or Möbius transformations) generalises the other two frameworks as demonstrated by Vuillemin [21]. Nielsen et al [18] showed that this framework can be used to represent quasi-normalised floating point [23].

We introduce here a new, feasible and incremental representation of the extended real numbers based on the composition of linear fractional transformations with *either all non-negative or all non-positive integer coefficients* [20].

Prototypes have been implemented in C++, Java and Miranda.

## 2 Linear Fractional Transformations

A natural way to represent a real number,  $r$  say, is by a sequence of nested rational intervals  $\{[p_n, q_n] : n \in \mathbb{N}\}$  enclosing  $r$  such that the sequence of interval lengths converges to zero [8, 17]:

$$[p_0, q_0] \supseteq [p_1, q_1] \supseteq [p_2, q_2] \supseteq [p_3, q_3] \supseteq \dots$$

Let  $\mathbb{R}$  denote the set of real numbers with the Euclidean topology,  $\mathbb{R}^\infty$  the one point compactification of  $\mathbb{R}$  and  $[0, \infty]$  the one point compactification of the non-negative real numbers  $[0, \infty)$ . The closed intervals  $[a, b]$  in  $\mathbb{R}^\infty$  are defined as the points from  $a$  to  $b$  in the numerically increasing direction, possibly including  $\infty$ .

Let us make the following convenient definitions:

$$\begin{aligned} \mathbb{V} &= \left\{ \begin{pmatrix} a \\ b \end{pmatrix} : a \neq 0 \text{ or } b \neq 0 \right\} \\ \mathbb{M} &= \left\{ \begin{pmatrix} a & c \\ b & d \end{pmatrix} : \begin{vmatrix} a & c \\ b & d \end{vmatrix} \neq 0 \right\} \\ \mathbb{T} &= \left\{ \begin{pmatrix} a & c & e & g \\ b & d & f & h \end{pmatrix} : \begin{vmatrix} ax + e & cx + g \\ bx + f & dx + h \end{vmatrix} \text{ and } \begin{vmatrix} ay + c & ey + g \\ by + d & fy + h \end{vmatrix} \text{ are non-trivial} \right\} \end{aligned}$$

Here  $\begin{vmatrix} a & c \\ b & d \end{vmatrix}$  refers to the determinant of the matrix  $\begin{pmatrix} a & c \\ b & d \end{pmatrix}$ , which is  $ad - bc$  of course.

**Definition 2.1** A 0-dimensional linear fractional transformation (lft) with real coefficients is a fraction in  $\mathbb{R}^\infty$ , namely a homogeneous coordinate representation of an extended real number,

$${}^t \begin{pmatrix} a \\ b \end{pmatrix} = \frac{a}{b} \quad (1)$$

where  $\begin{pmatrix} a \\ b \end{pmatrix} \in \mathbb{V}$ . A 1-dimensional lft with real coefficients is a function from  $\mathbb{R}^\infty$  to  $\mathbb{R}^\infty$  with the general form

$${}^t \begin{pmatrix} a & c \\ b & d \end{pmatrix} (x) = \frac{ax + c}{bx + d} \quad (2)$$

where  $\begin{pmatrix} a & c \\ b & d \end{pmatrix} \in \mathbb{M}$ . A 2-dimensional lft with real coefficients is a function from  $\mathbb{R}^\infty \times \mathbb{R}^\infty$  to  $\mathbb{R}^\infty$  with the general form

$${}^t \begin{pmatrix} a & c & e & g \\ b & d & f & h \end{pmatrix} (x, y) = \frac{axy + cx + ey + g}{bxy + dx + fy + h} \quad (3)$$

where  $\begin{pmatrix} a & c & e & g \\ b & d & f & h \end{pmatrix} \in \mathbb{T}$ .

In homogeneous coordinates, a 1-dimensional lft reduces to matrix multiplication

$$\begin{matrix} {}^t \begin{pmatrix} a & c \\ b & d \end{pmatrix} : & \mathbb{R}^\infty & \rightarrow & \mathbb{R}^\infty \\ & \begin{pmatrix} p \\ q \end{pmatrix} & \mapsto & \begin{pmatrix} ap + cq \\ bp + dq \end{pmatrix}. \end{matrix}$$

Therefore, it is convenient to drop the  $t$  in Equations (1), (2) and (3), We will also refer to the coefficients of a 0-dimensional lft as a vector, the coefficients of a 1-dimensional lft as a matrix and the coefficients of a 2-dimensional lft as a tensor.

**Definition 2.2** The information  $\text{Info}(P)$  contained by an arbitrary lft  $P$  is the interval in  $\mathbb{R}^\infty$  defined by  $\text{Info}(V) = \{V\}$ ,  $\text{Info}(M) = M([0, \infty])$  and  $\text{Info}(T) = T([0, \infty], [0, \infty])$ .

Consider a vector  $V$  as a pair of numbers denoted  $(V_0, V_1)$ ,

$$\begin{pmatrix} a \\ b \end{pmatrix} \equiv (a, b)$$

and consider a matrix  $M$  as a pair of vectors denoted  $(M_0, M_1)$ ,

$$\begin{aligned} \begin{pmatrix} a & c \\ b & d \end{pmatrix} &\equiv \left( \begin{pmatrix} a \\ b \end{pmatrix}, \begin{pmatrix} c \\ d \end{pmatrix} \right) \\ &\equiv ((a, b), (c, d)), \end{aligned}$$

and consider a tensor  $T$  as a pair of matrices denoted  $(T_0, T_1)$ ,

$$\begin{aligned} \begin{pmatrix} a & c & e & g \\ b & d & f & h \end{pmatrix} &\equiv \left( \begin{pmatrix} a & c \\ b & d \end{pmatrix}, \begin{pmatrix} e & g \\ f & h \end{pmatrix} \right) \\ &\equiv \left( \left( \begin{pmatrix} a \\ b \end{pmatrix}, \begin{pmatrix} c \\ d \end{pmatrix} \right), \left( \begin{pmatrix} e \\ f \end{pmatrix}, \begin{pmatrix} g \\ h \end{pmatrix} \right) \right) \\ &\equiv (((a, b), (c, d)), ((e, f), (g, h))). \end{aligned}$$

The first important observation is that  $\text{Info}(P)$  is computationally easy to evaluate. The endpoints of the interval  $\text{Info}(P)$  (provided it has any) correspond to a particular pair of the vectors that make up the coefficients. For a matrix  $M = \begin{pmatrix} a & c \\ b & d \end{pmatrix}$ , we have

$$\text{Info}(M) = \begin{cases} \left[ \frac{a}{b}, \frac{c}{d} \right] & \text{if } \det(M) < 0 \\ \left[ \frac{c}{d}, \frac{a}{b} \right] & \text{if } \det(M) > 0 \end{cases}$$

and for a tensor  $T = \begin{pmatrix} a & c & e & g \\ b & d & f & h \end{pmatrix}$ , we have

$$\text{Info}(T) = \text{Info} \left( \begin{pmatrix} a & c \\ b & d \end{pmatrix} \right) \cup \text{Info} \left( \begin{pmatrix} e & g \\ f & h \end{pmatrix} \right) \cup \text{Info} \left( \begin{pmatrix} a & e \\ b & f \end{pmatrix} \right) \cup \text{Info} \left( \begin{pmatrix} c & g \\ d & h \end{pmatrix} \right).$$

**Definition 2.3** *An arbitrary lft  $P$  satisfies the refinement property, denoted  $\mathcal{R}(P)$ , if the coefficients of  $P$  are all non-negative or all non-positive.*

Let us define  $\mathbb{V}^+ = \{V \in \mathbb{V} : \mathcal{R}(V)\}$ ,  $\mathbb{M}^+ = \{M \in \mathbb{M} : \mathcal{R}(M)\}$  and  $\mathbb{T}^+ = \{T \in \mathbb{T} : \mathcal{R}(T)\}$ .

The second important observation is that  $\text{Info}(P) \subseteq [0, \infty]$  if and only if  $P$  satisfies the refinement property.

Therefore, the composition  $Q \circ P$  of arbitrary lft's  $P$  and  $Q$  where  $P$  satisfies the refinement property corresponds to interval refinement. This is because, in this case,  $\text{Info}(Q \circ P) \subseteq \text{Info}(Q)$ . In other words,  $P$  refines the information given by  $Q$ .

### 3 Normal Products

**Proposition 3.1** *Given two rational intervals  $I = \text{Info}(M)$  and  $J = \text{Info}(N)$  represented by the two matrices  $M$  and  $N$ , then  $I \subseteq J$  if and only if there exists a matrix  $K$  with integer coefficients satisfying the refinement property such that  $M = NK$ .*

Therefore any real number can be represented as the intersection

$$\bigcap_{n \geq 0} M_0 M_1 M_2 \dots M_n ([0, \infty])$$

for a sequence of matrices  $M_n$  where  $M_n$  satisfies the refinement property for  $n \geq 1$ . We can denote this real number by an infinite product of matrices

$$\begin{pmatrix} a_0 & c_0 \\ b_0 & d_0 \end{pmatrix} : \begin{pmatrix} a_1 & c_1 \\ b_1 & d_1 \end{pmatrix} : \begin{pmatrix} a_2 & c_2 \\ b_2 & d_2 \end{pmatrix} : \dots \quad (4)$$

where  $M_n = \begin{pmatrix} a_n & c_n \\ b_n & d_n \end{pmatrix}$  and  $M_0 \in \mathbb{M}$  and  $M_n \in \mathbb{M}^+$  with  $n \geq 1$ . We will call this an *infinite normal product*. Notice however that an infinite normal product does not in general represent a point, although it always represents an interval. However, we are only interested in infinite normal products that converge to a point. We will call the first matrix  $M_0$  the *sign matrix* and the others matrices  $M_n$  with  $n \geq 1$  we will call the *digit matrices*. A singular matrix is in fact a constant that can be replaced by a vector, thus terminating the product; this we will call a *finite normal product*. Thus, a finite normal product represents a rational number, whereas an infinite normal product may represent any number.

The characteristic of a normal product is that the sign matrix identifies the interval in  $\mathbb{R}^\infty$  of interest, while the digit matrices refine this interval to a point. Therefore, we can consider two types called *signed normal products (snp)* and *unsigned normal products (unp)*. So, we can think of a signed normal product as a sign matrix followed by an unsigned normal product. Additionally, an unsigned normal product consists of a digit matrix followed by an unsigned normal products recursively. So, signed normal products represent real numbers in  $\mathbb{R}^\infty$ , while unsigned normal products represent real numbers in  $[0, \infty]$ . In BNF, we have

$$\begin{aligned} \text{snp}_0 & ::= \alpha \mid \beta : \text{unp}_{\phi(0,\beta)} \quad \text{where } \alpha \in \Phi_v(0) \text{ and } \beta \in \Phi_m(0) \\ \text{unp}_\tau & ::= \gamma \mid \delta : \text{unp}_{\phi(\tau,\delta)} \quad \text{where } \gamma \in \Phi_v(\tau) \text{ and } \delta \in \Phi_m(\tau) \text{ and } \tau \geq 1. \end{aligned}$$

We call  $\Phi_v(\tau)$  and  $\Phi_m(\tau)$  the *state sets*. They contain the next allowable vectors and matrices in the normal product given that we are in state  $\tau$ . The initial state is 0. Therefore  $\Phi_v(0)$  and  $\Phi_m(0)$  contain the allowable sign vectors and matrices.  $\phi(\tau, M)$  is called the *state transition function* and specifies the next state given that the previous state is  $\tau$  and the previous matrix in the normal product is  $M$ . Taken together, the

state sets and the state transition function restrict the allowable sequences of vectors and matrices in a normal product.

For the *most general normal product* we define the *state sets* to be

$$\begin{aligned}\Phi_v(0) &= \mathbf{V} \\ \Phi_v(1) &= \mathbf{V}^+ \\ \Phi_m(0) &= \mathbf{M} \\ \Phi_m(1) &= \mathbf{M}^+\end{aligned}$$

and the *state transition function*  $\phi(\tau, M)$  by

$\phi(\tau, M)$	$\tau = 0$	$\tau = 1$
$M \in \Phi_m(0)$	1	
$M \in \Phi_m(1)$		1

This notation may seem excessive at this point, but it will prove useful when we define various notions of *exact floating point* later.

## 4 Expression Trees

Let us generalise normal products even further to *expression trees* by including tensors. We can do this by constructing *signed expression trees (sexp)* and *unsigned expression trees (uexp)* according to the BNF prescription

$$\begin{aligned}\text{sexp} &::= \mathbf{V} \mid \mathbf{M}(\text{uexp}) \mid \mathbf{T}(\text{uexp}, \text{uexp}) \\ \text{uexp} &::= \mathbf{V}^+ \mid \mathbf{M}^+(\text{uexp}) \mid \mathbf{T}^+(\text{uexp}, \text{uexp}).\end{aligned}$$

If an expression  $E$  is a vector then it simply represents a rational number.

If an expression  $E$  is a matrix  $M$  applied to another expression  $E'$  then it represents an interval which is a subset of  $\text{Info}(M)$ . The information in the matrix  $M$  can be refined by incorporating information from the expression  $E'$  according to equations (5) and (6). We call this process *matrix absorption*.

If an expression  $E$  is a tensor  $T$  applied to two other expressions  $E'$  and  $E''$  then it represents an interval which is a subset of  $\text{Info}(T)$ . Again, the information in the tensor  $T$  can be refined by incorporating information from the expressions  $E'$  and  $E''$  according to the equations (7) to (10). We call this process *tensor absorption*.

What happens if we meet a tensor in a subexpression? Absorption equations could be devised to cater for this scenario. However, this would lead to higher dimensional lft's and increased complexity. A solution is to force subexpression tensors to output information in the form of matrices that reflect the information contained in them according to equation (11). We call this process *tensor emission*.

## 5 Absorption and Emission Equations

In order to simplify composition of lft's of various dimensions, we define the *dot product*, the *left product* and the *right product*, denoted respectively by  $\cdot$ ,  $\textcircled{L}$  and  $\textcircled{R}$  as follows:

$$\begin{aligned} (M \cdot V)_i &= \sum_{j=0,1} M_{ij} V_j \\ (M \cdot N) &= (M \cdot N_0, M \cdot N_1) \\ (M \cdot T) &= (M \cdot T_0, M \cdot T_1) \\ T \textcircled{R} V &= (T_0 \cdot V, T_1 \cdot V) \\ T \textcircled{R} M &= (T_0 \cdot M, T_1 \cdot M) \\ T \textcircled{L} V &= T^\top \textcircled{R} V \\ T \textcircled{L} M &= (T^\top \textcircled{R} M)^\top \end{aligned}$$

where  $T^\top$  indicates the transpose of  $T$  defined by swapping its middle two columns. Note that dot product is just conventional matrix multiplication. Let us also define the *mediant* of a matrix by

$$\overline{\begin{pmatrix} a & c \\ b & d \end{pmatrix}} = \begin{pmatrix} a+c \\ b+d \end{pmatrix}.$$

**Proposition 5.1** *The following matrix absorption equations hold:*

$$M(V) = M \cdot V \tag{5}$$

$$M(N(x)) = (M \cdot N)(x) \tag{6}$$

*The following tensor absorption equations hold:*

$$T(V, y) = \begin{cases} \overline{T \textcircled{L} V} & \text{if } |T \textcircled{L} V| = 0 \\ (T \textcircled{L} V)(y) & \text{if } |T \textcircled{L} V| \neq 0 \end{cases} \tag{7}$$

$$T(M(x), y) = (T \textcircled{L} M)(x, y) \tag{8}$$

$$T(x, V) = \begin{cases} \overline{T \textcircled{R} V} & \text{if } |T \textcircled{R} V| = 0 \\ (T \textcircled{R} V)(x) & \text{if } |T \textcircled{R} V| \neq 0 \end{cases} \tag{9}$$

$$T(x, M(y)) = (T \textcircled{R} M)(x, y). \tag{10}$$

**Proof** Here is the proof for one of the absorption equations:

$$\begin{aligned} T(x, V) &= \begin{pmatrix} a & c & e & g \\ b & d & f & h \end{pmatrix} \left( x, \begin{pmatrix} i \\ j \end{pmatrix} \right) \\ &= \frac{\frac{aix}{j} + cx + \frac{ei}{j} + g}{\frac{bix}{j} + dx + \frac{fi}{j} + h} \\ &= \frac{aix + cjx + ei + gj}{bix + djx + fi + hj} \\ &= \begin{pmatrix} ai + cj & ei + gj \\ bi + dj & fi + hj \end{pmatrix} (x) \end{aligned}$$

$$\begin{aligned}
&= \left( \left( \begin{pmatrix} a & c \\ b & d \end{pmatrix} \cdot \begin{pmatrix} i \\ j \end{pmatrix}, \begin{pmatrix} e & g \\ f & h \end{pmatrix} \cdot \begin{pmatrix} i \\ j \end{pmatrix} \right) (x) \\
&= (T_0 \cdot V, T_1 \cdot V)(x) \\
&= (T_{\mathbb{R}} \otimes V)(x) \\
&= \begin{cases} \overline{T_{\mathbb{R}} V} & \text{if } |T_{\mathbb{R}} V| = 0 \\ (T_{\mathbb{R}} V)(x) & \text{if } |T_{\mathbb{R}} V| \neq 0 \end{cases}
\end{aligned}$$

Note that the left and right products of a tensor with a vector may give a singular matrix, which is essentially a vector given by its median.  $\square$

The third important observation is that given a matrix  $M$  and an arbitrary lft  $P$ ,  $\text{Info}(M) \supseteq \text{Info}(P)$  if and only if  $\mathcal{R}(M^{-1} \cdot P)$  where  $M^{-1}$  is the matrix inverse of  $M$ . The scaling invariance of lft's means that we can simplify the definition of the matrix inverse to

$$\begin{pmatrix} a & c \\ b & d \end{pmatrix}^{-1} = \begin{pmatrix} d & -c \\ -b & a \end{pmatrix}.$$

**Proposition 5.2** *The following emission equations hold:*

$$\begin{aligned}
V &= E(E^{-1} \cdot V) && \text{if } \text{Info}(E) \supseteq \text{Info}(V) \\
M(x) &= E((E^{-1} \cdot M)(x)) && \text{if } \text{Info}(E) \supseteq \text{Info}(M) \\
T(x, y) &= E((E^{-1} \cdot T)(x, y)) && \text{if } \text{Info}(E) \supseteq \text{Info}(T)
\end{aligned} \tag{11}$$

*The conditions are necessary in order to ensure that  $(E^{-1} \cdot V)$ ,  $(E^{-1} \cdot M)$  and  $(E^{-1} \cdot T)$  satisfy the refinement property.*

Note that the fundamental aspects of these absorption and emission equations are essentially well known [18]. The novel aspects include the introduction of vectors and, very importantly, the repeated insistence on lft's satisfying the refinement property leading to a computationally efficient test for interval inclusion.

## 6 Tensor Absorption Strategy

For computing the value of  $T(x, y)$ , we need a **strategy** for deciding whether to absorb from  $x$  (*left absorption*) or from  $y$  (*right absorption*). All we know about  $x$  and  $y$  is that they are non-negative real numbers. So, what we need is a function  $\text{strategy}(T)$  which evaluates to **left**, **right** or **either**. By convention, we choose left absorption when we have a free choice. This enables algorithms to be made in the knowledge that there is a preferred absorption direction. See the  $\tan(x)$  example in section 12.



## 6.1 The Outcome Minimisation Strategy

One approach is to minimise the possible range of outcomes. The range of outcomes for a left absorption is given by

$$L(T) = \sup \left\{ d \left( \frac{ax + e}{bx + f}, \frac{cx + g}{dx + h} \right) : x \in [0, \infty] \right\}$$

and the range of outcomes for a right absorption is given by

$$R(T) = \sup \left\{ d \left( \frac{ay + c}{by + d}, \frac{ey + g}{fy + h} \right) : y \in [0, \infty] \right\}$$

where  $d(x, y)$  is a suitable metric such as

$$d(x, y) = \left| \frac{x}{x+1} - \frac{y}{y+1} \right|.$$

The exact formula involves a square root. However, a good approximation in practice for  $L(T)$  and  $R(T)$ , which is exact in many cases, is given by

$$\begin{aligned} L(T) &\approx \text{width}(J \cdot T^\top) \\ R(T) &\approx \text{width}(J \cdot T) \end{aligned}$$

where

$$\begin{aligned} J &= \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \\ \text{width}(T) &= \max \{ |T_0(0) - T_1(0)|, |T_0(\infty) - T_1(\infty)| \}. \end{aligned}$$

So, the strategy is:

$$\text{strategy}(T) = \begin{cases} \text{left} & \text{if } L(T) < R(T) \\ \text{either} & \text{if } L(T) = R(T) \\ \text{right} & \text{if } L(T) > R(T) \end{cases}$$

## 6.2 The Information Overlap Strategy

An alternative approach is to consider the pairs of matrices in tensor  $T$  and the pairs of matrices in tensor  $T^\top$ .

During left absorption the information in each of the matrices in tensor  $T^\top$  is refined, while during right absorption the information in each of the matrices in tensor  $T$  is refined.

Absorption is more effective if the pair of matrices under consideration have overlapping information because then the overall information refinement for the tensor is likely to be more dramatic on average. A more important observation is that absorption eventually leads to an empty intersection of information thus ensuring that the strategy is fair.

Let us define the function  $\text{overlap}(T)$  by

$$\begin{aligned} \text{overlap} : \mathbb{T} &\rightarrow \text{boolean} \\ T &\mapsto \text{Info}(T_0) \cap \text{Info}(T_1) \neq \emptyset \end{aligned}$$

Given  $T = \begin{pmatrix} a & c & e & g \\ b & d & f & h \end{pmatrix}$ , it can be shown that the boolean expression

$$\text{Info}(T_0) \cap \text{Info}(T_1) = \emptyset$$

is equivalent to

$$\text{Info} \left( \begin{pmatrix} d & -c \\ b & -a \end{pmatrix} \cdot \begin{pmatrix} e & g \\ f & h \end{pmatrix} \right) \subseteq (0, \infty).$$

It is always the case that at least one of  $\text{overlap}(T)$  and  $\text{overlap}(T^\top)$  is true. Therefore, a straightforward strategy is:

$$\text{strategy}(T) = \begin{cases} \text{left} & \text{if not } \text{overlap}(T) \\ \text{either} & \text{if } \text{overlap}(T) \text{ and } \text{overlap}(T^\top) \\ \text{right} & \text{if not } \text{overlap}(T^\top) \end{cases}$$

## 7 Unbiased Exact Floating Point

Storing a real number as a general normal product may be efficient in certain cases, but in general we need to ensure that the digit matrices in an infinite normal product refine the sign matrix at a steady rate. Also, we need a way to control the size of the integer coefficients in tensors. Both these objectives can be achieved by restricting the state sets.

In particular, we put forward the following important representation that we shall call *base  $B$  unbiased exact floating point*

$$S_\sigma : D_{d_1}^B : D_{d_2}^B : D_{d_3}^B : \dots$$

where

$$\begin{aligned} \sigma &\in \mathbb{S} = \{+, \infty, -, 0\} \\ d_i &\in \mathbb{D} = \{n - B, B - n : n \in \mathbb{N} \cap [1, B]\} \\ S_+ &= \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \\ S_\infty &= \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \\ S_- &= \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \end{aligned}$$

$$S_0 = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}$$

$$D_d^B = \begin{pmatrix} B+d+1 & B+d-1 \\ B-d-1 & B-d+1 \end{pmatrix}.$$

Therefore  $S_\sigma$  is the *sign matrix* while  $D_{d_i}^B$  for  $i \geq 1$  make up the *digit matrices* in base  $B$ . If the digit matrices begin with a block of  $n$   $D_{1-B}^B$ 's or  $D_{B-1}^B$ 's, we call this the *exponent* because it indicates the order of magnitude of the real number in question. The rest of the digit matrices is called the *mantissa*.

The sign matrices have been chosen such that they divide  $\mathbb{R}^\infty$  in a natural, unbiased and redundant manner;  $\text{Info}(S_+) = [0, \infty]$ ,  $\text{Info}(S_\infty) = [1, -1]$ ,  $\text{Info}(S_-) = [\infty, 0]$  and  $\text{Info}(S_0) = [-1, 1]$ . Recall that redundancy is important for computability. Also, they have positive determinants, which is convenient for implementation, and form a cyclic group (up to scaling invariance)

$$\begin{aligned} S_\infty^0 &= S_+ \\ S_\infty^1 &= S_\infty \\ S_\infty^2 &= S_- \\ S_\infty^3 &= S_0 \\ S_\infty^4 &= S_\infty^0, \end{aligned}$$

which is mathematically useful for deriving algorithms.

The digit matrices satisfy the following property, which is reminiscent of redundant sequences of digits proposed by Avizienis [1].

$$D_{d_1}^B \cdot D_{d_2}^B \cdot \dots \cdot D_{d_n}^B = \begin{pmatrix} B^n + m + 1 & B^n + m - 1 \\ B^n - m - 1 & B^n - m + 1 \end{pmatrix}$$

where

$$m = \sum_{i=1}^n d_i B^{n-i}.$$

In fact, the sign matrix  $S_0$  followed by  $D_{d_1}^B : D_{d_2}^B : D_{d_3}^B : \dots$  corresponds exactly to redundant sequences of binary digits as pointed out by Nielsen et al [18]. The novelty here is that we extend the representation to any base  $B$  and to the other three signs;  $+$ ,  $\infty$  and  $-$ . For  $S_0$ , we have

$$\begin{aligned} S_0 D_{d_1}^B \cdot D_{d_2}^B \cdot \dots \cdot D_{d_n}^B &= \begin{pmatrix} m+1 & m-1 \\ B^n & B^n \end{pmatrix} \\ \text{Info}(S_0 D_{d_1}^B \cdot D_{d_2}^B \cdot \dots \cdot D_{d_n}^B) &= \left[ \frac{m-1}{B^n}, \frac{m+1}{B^n} \right]. \end{aligned} \quad (12)$$

The origin of  $D_d^B$  becomes clearer when we consider the following commutative diagram:

$$\begin{array}{ccc} [0, \infty] & \xrightarrow{D_d^B} & [0, \infty] \\ \downarrow S_0 & & \downarrow S_0 \\ [-1, 1] & \xrightarrow{R_d^B} & [-1, 1] \end{array}$$

where

$$\begin{aligned} R_d^B &= S_0 \cdot D_d^B \cdot S_0^{-1} \\ &= \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} B+d+1 & B+d-1 \\ B-d-1 & B-d+1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \\ &= \begin{pmatrix} 1 & d \\ 0 & B \end{pmatrix} \\ R_d^B(x) &= \frac{x+d}{B} \end{aligned}$$

So,  $R_d^B$  are the standard affine maps corresponding to the digits in the representation of the real numbers over the interval  $[-1, 1]$  by redundant sequences of digits. In particular,  $R_{-1}^2$  maps  $[-1, 1]$  to the left half  $[-1, 0]$ ,  $R_0^2$  maps  $[-1, 1]$  to the middle half  $[-\frac{1}{2}, \frac{1}{2}]$  and  $R_1^2$  maps  $[-1, 1]$  to the right half  $[0, 1]$ .

So, the state sets for unbiased exact floating point are

$$\begin{aligned} \Phi_v(0) &= \emptyset \\ \Phi_v(1) &= \emptyset \\ \Phi_m(0) &= \{S_\sigma : \sigma \in \mathbb{S}\} \\ \Phi_m(1) &= \{D_n^B : n \in \mathbb{D}\} \end{aligned}$$

and the state transition function  $\phi(\tau, M)$  is given by

$\phi(\tau, M)$	$\tau = 0$	$\tau = 1$
$M \in \Phi_m(0)$	1	
$M \in \Phi_m(1)$		1

A useful hybrid between general normal products and unbiased exact floating point involves setting  $\Phi_v(0)$  and  $\Phi_v(1)$  to

$$\begin{aligned} \Phi_v(0) &= \left\{ \begin{pmatrix} a \\ b \end{pmatrix} : \begin{pmatrix} a \\ b \end{pmatrix} \in \mathbb{V} \text{ and } |a|, |b| < \Omega \right\} \\ \Phi_v(1) &= \left\{ \begin{pmatrix} a \\ b \end{pmatrix} : \begin{pmatrix} a \\ b \end{pmatrix} \in \mathbb{V}^+ \text{ and } |a|, |b| < \Omega \right\} \end{aligned}$$

where  $\Omega$  may correspond to the maximum integer that fits into a computer word.

Interesting bases include 2 with digits  $\{-1, 0, 1\}$  and the golden ratio  $\phi = \frac{1+\sqrt{5}}{2}$  with digits  $\{1 - \phi, \phi - 1\}$  [19, 6].

This representation of the real numbers leads to efficient algorithms for a wide range of operations including the transcendental functions. The only exception is for addition and subtraction. In which case, we need biased exact floating point.

## 8 Biased Exact Floating Point

Biased exact floating point leads to efficient algorithms for addition and subtraction that are essentially the same as those for redundant sequences of digits with the added benefit of incrementality even at the exponent level. This means that infinity can be handled as if it were any other real number.

We restrict the allowable sequences of matrices as follows

$$\left( \begin{array}{l} S_+ : \left\{ \begin{array}{l} L \\ H : X^* : F_n \end{array} \right\} \\ S_\infty : Z^* : \left\{ \begin{array}{l} I : X^* : F_n \\ \hat{I} : \hat{X}^* : \hat{F}_n \end{array} \right\} \\ S_- : \left\{ \begin{array}{l} \hat{L} \\ \hat{H} : \hat{X}^* : \hat{F}_n \end{array} \right\} \\ S_0 \end{array} \right) : \left\{ \begin{array}{l} D_{1-B}^B \\ \vdots \\ D_{B-1}^B \end{array} \right\}^*$$

where  $M^*$  indicates zero or more occurrences of  $M$  and  $\left\{ \begin{array}{l} P \\ Q \end{array} \right\}$  indicates a choice between the sequences  $P$  and  $Q$  and

$$\begin{aligned} n &\in \mathbb{F} = \mathbb{N} \cap [1, B) \\ L &= \begin{pmatrix} 2 & 0 \\ 1 & 1 \end{pmatrix} \\ H &= \begin{pmatrix} B & 1 \\ 0 & 1 \end{pmatrix} \\ X &= \begin{pmatrix} B^2 & B-1 \\ 0 & B \end{pmatrix} \\ F_n &= \begin{pmatrix} n+1 & n-1 \\ B & B \end{pmatrix} \\ I &= \begin{pmatrix} B & 0 \\ B & 2 \end{pmatrix} \\ Z &= \begin{pmatrix} B+1 & B-1 \\ B-1 & B+1 \end{pmatrix} \end{aligned}$$

$$\widehat{\begin{pmatrix} a & c \\ b & d \end{pmatrix}} = \begin{pmatrix} d & b \\ c & a \end{pmatrix}.$$

These matrices are such that once we reach the matrices  $D_{1-B}^B$  to  $D_{B-1}^B$  the nested intervals correspond to conventional  $B$ -adic intervals (as in equation (12)). The following properties can be easily derived:

$$\begin{aligned} \text{Info}(S_+L) &= [0, 2] \\ \text{Info}(S_+HX^e) &= [B^e, \infty] \\ \text{Info}(S_+HX^eF_n) &= [nB^e, (n+2)B^e] \\ \text{Info}(S_\infty Z^e) &= [B^e, -B^e] \\ \text{Info}(S_\infty Z^e I) &= [B^e, \infty] \\ \text{Info}(S_\infty Z^e \hat{I}) &= [\infty, -B^e] \\ \text{Info}(S_- \hat{L}) &= [-2, 0] \\ \text{Info}(S_- \hat{H} \hat{X}^e) &= [\infty, -B^e] \\ \text{Info}(S_- \hat{H} \hat{X}^e \hat{F}_n) &= [-(n+2)B^e, -nB^e]. \end{aligned}$$

It should be pointed out that Nielsen et al [18] considered a similar idea for quasi-normalised floating point. However, we restrict the digit matrices to ones satisfying the refinement property. This means that our algorithms for interval inclusion are much more efficient. Also, our representation includes three types of infinity; namely positive infinity, negative infinity and an infinity with an unknown sign. This means that an algorithm for addition can emit matrices without necessarily first absorbing the entire exponent part of the two arguments.

To summarise, the state sets are

$$\begin{aligned} \Phi_v(n) &= \emptyset \\ \Phi_m(0) &= \{S_\sigma : \sigma \in \mathbb{S}\} \\ \Phi_m(1) &= \{D_n^B : n \in \mathbb{D}\} \\ \Phi_m(2) &= \{L, H\} \\ \Phi_m(3) &= \{I, Z, \hat{I}\} \\ \Phi_m(4) &= \{\hat{L}, \hat{H}\} \\ \Phi_m(5) &= \{X\} \cup \{F_n : n \in \mathbb{F}\} \\ \Phi_m(6) &= \{\hat{X}\} \cup \{\hat{F}_n : n \in \mathbb{F}\}. \end{aligned}$$

and the state transition function  $\phi(\tau, M)$  is given by

$\phi(\tau, M)$	$\tau = 0$	$\tau = 1$	$\tau = 2$	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 6$
$M = S_+$	2						
$M = S_\infty$	3						
$M = S_-$	4						
$M = S_0$	1						
$M = D_n^B$		1					
$M = L$			1				
$M = H$			5				
$M = I$				5			
$M = Z$				3			
$M = \hat{I}$				6			
$M = \hat{L}$					1		
$M = \hat{H}$					6		
$M = X$						5	
$M = F_n$						1	
$M = \hat{X}$							6
$M = \hat{F}_n$							1

## 9 Normalisation Algorithms

In this section, we describe a basic set of lazy functions that allow signed normal products to be generated from

(0-dimensional lft) rational numbers,

(1-dimensional lft) matrices applied to signed general normal products and

(2-dimensional lft) tensors applied to signed general normal products.

The signed normal products generated follow the restrictions specified by the state sets  $\Phi_v(\tau)$  and  $\Phi_m(\tau)$  for  $\tau \in \mathbb{N}$  and the state transition function  $\phi(\tau, M)$  for  $\tau \in \mathbb{N}$  and  $M \in \mathbb{M}$ .

### 9.1 Vector Algorithm

The expression  $\text{svector}_0(V)$  evaluates the vector  $V \in \mathbb{V}$  to a signed normal product.

$$\begin{aligned}
 \text{svector}_0 : \mathbb{V} &\rightarrow \text{snp}_0 \\
 V &\mapsto V \text{ if } V \in \Phi_v(0) \\
 V &\mapsto S : \text{uvector}_{\phi(0,S)}(S^{-1} \cdot V) \\
 &\quad \text{if } \exists S \in \Phi_m(0) \text{ such that } \text{Info}(S) \supseteq \text{Info}(V)
 \end{aligned}$$

In words, if  $V$  is an allowable sign vector then it is left unchanged otherwise find a sign matrix  $S$  that can be emitted and continue with the remainder.

The expression  $\text{uvector}_\tau(V)$  evaluates the vector  $V \in \mathbb{V}^+$  to a state  $\tau$  unsigned normal product.

$$\begin{aligned} \text{uvector}_\tau : \mathbb{V}^+ &\rightarrow \text{unp}_\tau \\ V &\mapsto V \text{ if } V \in \Phi_v(\tau) \\ V &\mapsto D : \text{uvector}_{\phi(\tau,D)}(D^{-1} \cdot V) \\ &\quad \text{if } \exists D \in \Phi_m(\tau) \text{ such that } \text{Info}(D) \supseteq \text{Info}(V) \end{aligned}$$

In words, if  $V$  is an allowable digit vector in state  $\tau$  then leave it unchanged otherwise find a digit matrix  $D$  that can be emitted in this state and continue with the remainder.

## 9.2 Matrix Algorithm

The expression  $\text{smatrix}_s(M, E_x)$  evaluates the matrix  $M \in \mathbb{M}$  and the signed general normal product  $E_x$  to a signed normal product.

$$\begin{aligned} \text{smatrix}_s : \mathbb{M} \times \text{sexp} &\rightarrow \text{snp}_0 \\ (M, V_x) &\mapsto \text{svector}_0(M \cdot V_x) \\ (M, S_x : U_x) &\mapsto \text{smatrixu}_0(M \cdot S_x, U_x) \end{aligned}$$

In words, the signed part of the signed general normal product  $E_x$  must be absorbed into  $M$  before any further steps can be taken.

The expression  $\text{smatrixu}_0(M, E_x)$  evaluates the matrix  $M \in \mathbb{M}$  and the unsigned general normal product  $E_x$  to a signed normal product.

$$\begin{aligned} \text{smatrixu}_0 : \mathbb{M} \times \text{uexp} &\rightarrow \text{snp}_0 \\ (M, V_x) &\mapsto M : \text{uvector}_{\phi(0,M)}(V_x) \text{ if } M \in \Phi_m(0) \\ (M, D_x : U_x) &\mapsto M : \text{umatrixu}_{\phi(0,M)}(D_x, U_x) \text{ if } M \in \Phi_m(0) \\ (M, E_x) &\mapsto S : \text{umatrixu}_{\phi(0,S)}(S^{-1} \cdot M, E_x) \\ &\quad \text{if } \exists S \in \Phi_m(0) \text{ such that } \text{Info}(S) \supseteq \text{Info}(M) \\ (M, V_x) &\mapsto \text{svector}_0(M \cdot V_x) \\ (M, D_x : U_x) &\mapsto \text{smatrixu}_0(M \cdot D_x, U_x) \end{aligned}$$

In words, if  $M$  is an allowable sign matrix then it is left unchanged otherwise try to find a sign matrix  $S$  that can be emitted and continue with the remainder. Otherwise, more information must be absorbed from  $E_x$ .

The expression  $\text{umatrixu}_\tau(M, E_x)$  evaluates the matrix  $M \in \mathbb{M}^+$  and the unsigned general normal product  $E_x$  to a state  $\tau$  unsigned normal product.

$$\text{umatrixu}_\tau : \mathbb{M}^+ \times \text{uexp} \rightarrow \text{unp}_\tau$$



$$\begin{aligned}
(M, V_x) &\mapsto M : \text{uvector}_{\phi(\tau, M)}(V_x) \text{ if } M \in \Phi_m(\tau) \\
(M, D_x : U_x) &\mapsto M : \text{umatrixu}_{\phi(\tau, M)}(D_x, U_x) \text{ if } M \in \Phi_m(\tau) \\
(M, E_x) &\mapsto D : \text{umatrixu}_{\phi(\tau, D)}(D^{-1} \cdot M, E_x) \\
&\quad \text{if } \exists D \in \Phi_m(\tau) \text{ such that } \text{Info}(D) \supseteq \text{Info}(M) \\
(M, V_x) &\mapsto \text{uvector}_{\tau}(M \cdot V_x) \\
(M, D_x : U_x) &\mapsto \text{umatrixu}_{\tau}(M \cdot D_x, U_x)
\end{aligned}$$

In words, if  $M$  is an allowable digit matrix in state  $\tau$  then it is left unchanged otherwise try to find a digit matrix  $D$  that can be emitted in this state and continue with the remainder. Otherwise, more information must be absorbed from  $E_x$ .

### 9.3 Tensor Algorithm

The expression  $\text{stensors}_0(T, E_x, E_y)$  evaluates the tensor  $T \in \mathbb{T}$  and the signed general normal products  $E_x$  and  $E_y$  to a signed normal product.

$$\begin{aligned}
\text{stensors}_0 : \mathbb{T} \times \text{sexp} \times \text{sexp} &\rightarrow \text{snp}_0 \\
(T, V_x, V_y) &\mapsto \begin{cases} \text{svector}_0(\overline{T^{\mathbb{L}} V_x}) & \text{if } |T^{\mathbb{L}} V_x| = 0 \\ \text{svector}_0(T^{\mathbb{L}} V_x \cdot V_y) & \text{if } |T^{\mathbb{L}} V_x| \neq 0 \end{cases} \\
(T, V_x, S_y : U_y) &\mapsto \begin{cases} \text{svector}_0(\overline{T^{\mathbb{L}} V_x}) & \text{if } |T^{\mathbb{L}} V_x| = 0 \\ \text{smatrixu}_0(T^{\mathbb{L}} V_x \cdot S_y, U_y) & \text{if } |T^{\mathbb{L}} V_x| \neq 0 \end{cases} \\
(T, S_x : U_x, V_y) &\mapsto \begin{cases} \text{svector}_0(\overline{T^{\mathbb{R}} V_y}) & \text{if } |T^{\mathbb{R}} V_y| = 0 \\ \text{smatrixu}_0(T^{\mathbb{R}} V_y \cdot S_x, U_x) & \text{if } |T^{\mathbb{R}} V_y| \neq 0 \end{cases} \\
(T, S_x : U_x, S_y : U_y) &\mapsto \text{stensoru}_0(T^{\mathbb{L}} S_x^{\mathbb{R}} S_y, U_x, U_y)
\end{aligned}$$

In words, the signed part of the signed general normal products  $E_x$  and  $E_y$  must be absorbed into  $T$  before any further steps can be taken.

The expression  $\text{stensoru}_0(T, E_x, E_y)$  evaluates the tensor  $T \in \mathbb{T}$  and the unsigned general normal products  $E_x$  and  $E_y$  to a signed normal product.

$$\begin{aligned}
\text{stensoru}_0 : \mathbb{T} \times \text{uexp} \times \text{uexp} &\rightarrow \text{snp}_0 \\
(T, E_x, E_y) &\mapsto S : \text{utensoru}_{\phi(0, S)}(S^{-1} \cdot T, E_x, E_y) \\
&\quad \text{if } \exists S \in \Phi_m(0) \text{ such that } \text{Info}(S) \supseteq \text{Info}(T)
\end{aligned}$$

if  $\text{strategy}(T) \neq \text{right}$  then

$$\begin{aligned}
(T, V_x, U_y) &\mapsto \begin{cases} \text{svector}_0(\overline{T^{\mathbb{L}} V_x}) & \text{if } |T^{\mathbb{L}} V_x| = 0 \\ \text{smatrixu}_0(T^{\mathbb{L}} V_x, U_y) & \text{if } |T^{\mathbb{L}} V_x| \neq 0 \end{cases} \\
(T, D_x : U_x, U_y) &\mapsto \text{stensoru}_0(T^{\mathbb{L}} D_x, U_x, U_y)
\end{aligned}$$

if  $\text{strategy}(T) = \text{right}$  then

$$\begin{aligned}
(T, U_x, V_y) &\mapsto \begin{cases} \text{svector}_0(\overline{T^{\mathbb{R}} V_y}) & \text{if } |T^{\mathbb{R}} V_y| = 0 \\ \text{smatrixu}_0(T^{\mathbb{R}} V_y, U_x) & \text{if } |T^{\mathbb{R}} V_y| \neq 0 \end{cases} \\
(T, U_x, D_y : U_y) &\mapsto \text{stensoru}_0(T^{\mathbb{R}} D_y, U_x, U_y)
\end{aligned}$$

In words, try to find a sign matrix  $S$  that can be emitted and continue with the remainder. Otherwise, more information must be absorbed from  $E_x$  or  $E_y$  according to the value of  $\text{strategy}(T)$ . If  $\text{strategy}(T) = \mathbf{right}$  then absorb from  $E_y$  otherwise absorb from  $E_x$ .

The expression  $\text{utensoru}_\tau(T, E_x, E_y)$  evaluates the tensor  $T \in \mathbb{T}^+$  and the unsigned general normal products  $E_x$  and  $E_y$  to a state  $\tau$  unsigned normal product.

$$\begin{aligned} \text{utensoru}_\tau : \mathbb{T}^+ \times \text{uexp} \times \text{uexp} &\rightarrow \text{unp}_\tau \\ (T, E_x, E_y) &\mapsto D : \text{utensoru}_{\phi(\tau, D)}(D^{-1} \cdot T, E_x, E_y) \\ &\quad \text{if } \exists D \in \Phi_m(x) \text{ such that } \text{Info}(D) \supseteq \text{Info}(T) \end{aligned}$$

if  $\text{strategy}(T) \neq \mathbf{right}$  then

$$\begin{aligned} (T, V_x, U_y) &\mapsto \begin{cases} \text{uvector}_\tau(\overline{T \circledast V_x}) & \text{if } |T \circledast V_x| = 0 \\ \text{umatrixu}_\tau(T \circledast V_x, U_y) & \text{if } |T \circledast V_x| \neq 0 \end{cases} \\ (T, D_x : U_x, U_y) &\mapsto \text{utensoru}_\tau(T \circledast D_x, U_x, U_y) \end{aligned}$$

if  $\text{strategy}(T) = \mathbf{right}$  then

$$\begin{aligned} (T, U_x, V_y) &\mapsto \begin{cases} \text{uvector}_\tau(\overline{T \circledast V_y}) & \text{if } |T \circledast V_y| = 0 \\ \text{umatrixu}_\tau(T \circledast V_y, U_x) & \text{if } |T \circledast V_y| \neq 0 \end{cases} \\ (T, U_x, D_y : U_y) &\mapsto \text{utensoru}_\tau(T \circledast D_y, U_x, U_y) \end{aligned}$$

In words, try to find a digit matrix  $D$  that can be emitted in state  $\tau$  and continue with the remainder. Otherwise, more information must be absorbed from  $E_x$  or  $E_y$  according to the value of  $\text{strategy}(T)$ . If  $\text{strategy}(T) = \mathbf{right}$  then absorb from  $E_y$  otherwise absorb from  $E_x$ .

## 10 Flatten Algorithm

The basic set of functions above have to be used with care particularly where recursion is involved.

For example, consider  $\sqrt{2}$ . If we define the `sqrt2` recursively by

$$\text{sqrt2} = \text{smatrixs}_0 \left( \left( \begin{pmatrix} 1 & 2 \\ 1 & 1 \end{pmatrix}, \text{sqrt2} \right) \right)$$

then it will never output anything and never terminate.

If however, we define the `sqrt2` recursively by

$$\text{sqrt2} = \text{smatrixu}_0(S_+, x) \text{ where } x = \begin{pmatrix} 1 & 2 \\ 1 & 1 \end{pmatrix} : x$$

then it evaluates the answer as expected.

In general, given a mathematical formula we must be able to construct a signed expression tree and evaluate it in the following way using the function `sflatten`.

$$\text{sflatten} : \text{sexp} \rightarrow \text{snp}_0$$

$$\begin{aligned}
V &\mapsto \text{svector}_0(V) \\
M(E_x) &\mapsto \text{smatrixu}_0(M, \text{uflatten}(E_x)) \\
T(E_x, E_y) &\mapsto \text{stensoru}_0(T, \text{uflatten}(E_x), \text{uflatten}(E_y)) \\
\text{uflatten} : \text{uexp} &\rightarrow \text{uexp} \\
V &\mapsto V \\
M(E_x) &\mapsto M(\text{uflatten}(E_x)) \\
T(E_x, E_y) &\mapsto \text{utensoru}_1(T, \text{uflatten}(E_x), \text{uflatten}(E_y)).
\end{aligned}$$

Note that state 1 in  $\text{utensoru}_1$  can be replaced by any valid non-zero state. In words, the function  $\text{uflatten}$  merely flattens an unsigned expression tree into a general normal product while the function  $\text{sflatten}$  flattens a signed expression tree into a more specific normal product. The specific normal product is controlled by the state sets  $\Phi_v(\tau)$  and  $\Phi_m(\tau)$  for  $\tau \in \mathbb{N}$  and the state transition function  $\phi(\tau, M)$  for  $\tau \in \mathbb{N}$  and  $M \in \mathbb{M}$ .

So, for the simple example  $\sqrt{2}$ , we can define the expression tree  $\text{sqrt2}$  recursively by

$$\text{sqrt2} = \begin{pmatrix} 1 & 2 \\ 1 & 1 \end{pmatrix} (\text{sqrt2})$$

and then  $\text{sflatten}(\text{sqrt2})$  will convert it to the required normal product. For the general normal product, it remains unchanged. For unbiased exact floating point in base 2, the output is not unique due to redundancy in the representation, but it may look something like

$$S_+ : D_0^2 : D_1^2 : D_{-1}^2 : D_1^2 : D_0^2 : D_{-1}^2 : D_0^2 : D_0^2 : D_0^2 : D_0^2 : D_{-1}^2 : \dots$$

which represents the interval  $\left[\frac{1199}{849}, \frac{1200}{848}\right]$  thus far. For biased exact floating point in base 2, the output is also not unique due to redundancy in the representation, but it may look something like

$$S_+ : L : D_1^2 : D_{-1}^2 : D_1^2 : D_0^2 : D_1^2 : D_0^2 : D_1^2 : D_0^2 : D_0^2 : D_0^2 : \dots$$

which represents the interval  $\left[\frac{1447}{1024}, \frac{1449}{1024}\right]$  thus far.

## 11 Basic Arithmetic Operations

As already mentioned, we can consider the application of a 1-dimensional lft to a normal product as an arithmetic operation:

$$\begin{aligned}
M_{\text{neg}}(x) &= \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix} (x) = -x \\
M_{\text{rec}}(x) &= \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} (x) = \frac{1}{x}
\end{aligned}$$

Given a general 1-dimensional lft  $M(x)$ , we can simply use the function  $\text{smatrix}_{s_0}(M, x)$  to evaluate the result. So, negation and reciprocal can be defined by

$$\begin{aligned} \text{sneg} : \text{sexp} &\rightarrow \text{snp}_0 \\ x &\mapsto \text{smatrix}_{s_0}(M_{\text{neg}}, x) \\ \text{srec} : \text{sexp} &\rightarrow \text{snp}_0 \\ x &\mapsto \text{smatrix}_{s_0}(M_{\text{rec}}, x) \end{aligned}$$

Gosper [7] devised algorithms for the elementary arithmetic operations on continued fractions [22] using 2-dimensional lft's. The four most basic arithmetic operations can be represented as follows:

$$\begin{aligned} T_{\text{add}}(x, y) &= \begin{pmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} (x, y) = x + y \\ T_{\text{sub}}(x, y) &= \begin{pmatrix} 0 & 1 & -1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} (x, y) = x - y \\ T_{\text{mul}}(x, y) &= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} (x, y) = x \times y \\ T_{\text{div}}(x, y) &= \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} (x, y) = x \div y \end{aligned}$$

So, these operations can be defined by

$$\begin{aligned} \text{sadd} : \text{sexp} \times \text{sexp} &\rightarrow \text{snp}_0 \\ (x, y) &\mapsto \text{stensors}_0(T_{\text{add}}, x, y) \\ \text{ssub} : \text{sexp} \times \text{sexp} &\rightarrow \text{snp}_0 \\ (x, y) &\mapsto \text{stensors}_0(T_{\text{sub}}, x, y) \\ \text{smul} : \text{sexp} \times \text{sexp} &\rightarrow \text{snp}_0 \\ (x, y) &\mapsto \text{stensors}_0(T_{\text{mul}}, x, y) \\ \text{sdiv} : \text{sexp} \times \text{sexp} &\rightarrow \text{snp}_0 \\ (x, y) &\mapsto \text{stensors}_0(T_{\text{div}}, x, y) \end{aligned}$$

## 12 Continued Fractions

The development

$$a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \dots}}} \tag{13}$$

is called a *simple continued fraction* [3, 11].

The quantity

$$r_n = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \cdots + \frac{1}{a_n}}} \quad (14)$$

is called the  $n^{\text{th}}$  *approximant*. The  $0^{\text{th}}$  approximant is  $a_0$ . If the sequence  $r_n$  converges to a real number  $r$  then the continued fraction is said to be convergent and represent the number  $r$ .

Using the lft's

$$M_n(x) = \begin{pmatrix} a_n & 1 \\ 1 & 0 \end{pmatrix} (x) = a_n + \frac{1}{x} \quad (15)$$

we can generate the continued fraction in equation (13). Therefore, a continued fraction with non-negative coefficients for  $n \geq 1$  corresponds to the general normal product

$$M_0 : M_1 : M_2 : M_3 : \dots$$

Observe that  $(M_0 \cdot M_1 \cdot M_2 \cdot \dots \cdot M_n)([1, -1]) = (N_0 \cdot N_1 \cdot N_2 \cdot \dots \cdot N_n)([0, \infty])$  where

$$N_{n=0} = M_0 \cdot S_\infty = \begin{pmatrix} a_0 - 1 & a_0 + 1 \\ 1 & 1 \end{pmatrix}$$

$$N_{n \geq 1} = S_\infty^{-1} \cdot M_n \cdot S_\infty = \begin{pmatrix} a_n - 2 & a_n \\ a_n & a_n + 2 \end{pmatrix}.$$

We can see this more clearly in the following commutative diagram:

$$\begin{array}{ccccccc} [1, -1] & \xleftarrow{M_0} & [1, -1] & \xleftarrow{M_1} & [1, -1] & \xleftarrow{M_2} & [1, -1] \cdots \\ & & \uparrow S_\infty & & \uparrow S_\infty & & \uparrow S_\infty \\ & & [0, \infty] & \xleftarrow{N_1} & [0, \infty] & \xleftarrow{N_2} & [0, \infty] \cdots \\ & \swarrow N_0 & & & & & \end{array}$$

Therefore, a continued fraction with  $|a_n| \geq 2$  for  $n \geq 1$  corresponds to the general normal product

$$N_0 : N_1 : N_2 : N_3 : \dots$$

For example, a known formula for  $\tan(x)$  equates to

$$a_{n=0} = 0$$

$$a_{n \geq 1} = (-1)^{n+1} \frac{2n-1}{x}.$$

This leads to

$$N_{n=0} = \begin{pmatrix} -1 & 1 \\ 1 & 1 \end{pmatrix}$$

$$N_{n \geq 1} = \begin{pmatrix} 2n-1 + (-1)^n 2x & 2n-1 \\ 2n-1 & 2n-1 - (-1)^n 2x \end{pmatrix}.$$

Therefore, we can derive the following expression tree for  $\tan(S_0 : x)$

$$\begin{aligned} \tan & : \text{sexp} \rightarrow \text{sexp} \\ S_0 : x & \mapsto \begin{pmatrix} 1 & 1 & -1 & -1 \\ 0 & 2 & 2 & 0 \end{pmatrix} (x, \text{iterate}(\text{iterator}, 0, x)) \\ \\ \text{iterate} & : (\mathbb{N} \rightarrow \mathbb{T}^+) \times \mathbb{N} \times \text{uexp} \rightarrow \text{uexp} \\ & (T, n, x) \mapsto T(n)(x, \text{iterate}(T, n+1, x)) \\ \\ \text{iterator} & : \mathbb{N} \rightarrow \mathbb{T}^+ \\ n & \mapsto \begin{cases} \begin{pmatrix} 2n+5 & 2n+3 & 2n+1 & 2n+3 \\ 2n+3 & 2n+1 & 2n+3 & 2n+5 \end{pmatrix} & \text{if } n \text{ even} \\ \begin{pmatrix} 2n+1 & 2n+3 & 2n+5 & 2n+3 \\ 2n+3 & 2n+5 & 2n+3 & 2n+1 \end{pmatrix} & \text{if } n \text{ odd} \end{cases} \end{aligned}$$

Observe how the recursion has been defined in order to take advantage of the preference for left absorption in the tensor absorption strategies.

For example, consider  $\tan(\tan(1/3))$ . Represent  $1/3$  by  $\begin{pmatrix} 1 \\ 3 \end{pmatrix}$ . We first emit a sign matrix such as  $S_0 : \begin{pmatrix} 2 \\ 1 \end{pmatrix}$ . Therefore  $\tan(1/3)$  reduces to

$$\begin{pmatrix} 1 & 1 & -1 & -1 \\ 0 & 2 & 2 & 0 \end{pmatrix} \left( \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \text{iterate} \left( \text{iterator}, 0, \begin{pmatrix} 2 \\ 1 \end{pmatrix} \right) \right).$$

This leads to

$$\begin{pmatrix} 1 & 1 \\ 2 & 4 \end{pmatrix} \left( \text{iterate} \left( \text{iterator}, 0, \begin{pmatrix} 2 \\ 1 \end{pmatrix} \right) \right).$$

via left absorption because **strategy**  $\left( \begin{pmatrix} 1 & 1 & -1 & -1 \\ 0 & 2 & 2 & 0 \end{pmatrix} \right) = \mathbf{either}$  and left absorption is the preferred direction. Again we need a sign matrix and so this reduces to

$$S_0 : \begin{pmatrix} 3 & 5 \\ 1 & 3 \end{pmatrix} \left( \text{iterate} \left( \text{iterator}, 0, \begin{pmatrix} 2 \\ 1 \end{pmatrix} \right) \right).$$

Hence  $\tan(\tan(1/3))$  reduces to the signed expression tree

$$\begin{pmatrix} 1 & 1 & -1 & -1 \\ 0 & 2 & 2 & 0 \end{pmatrix} (x, \text{iterate}(\text{iterator}, 0, x))$$

where

$$x = \begin{pmatrix} 3 & 5 \\ 1 & 3 \end{pmatrix} \left( \text{iterate} \left( \text{iterator}, 0, \begin{pmatrix} 2 \\ 1 \end{pmatrix} \right) \right).$$

Finally, we apply the function `sflatten` for biased exact floating point in base 2 to this signed expression tree. The output is not unique due to redundancy in the representation, but it may look something like

$$S_+ : L_+ : D_{-1}^2 : D_{-1}^2 : D_1^2 : D_0^2 : D_0^2 : D_{-1}^2 : D_0^2 : D_0^2 : D_1^2 : D_{-1}^2 : \dots$$

which represents the interval  $\left[ \frac{369}{1024}, \frac{370}{1024} \right]$  thus far.

A survey of various other mathematical functions in the form of general normal products has been made by Potts [20].

## 13 Conclusion

In conclusion, may we draw attention to the novel ideas expressed in this document.

- The use of vectors to represent rational number finitely in this setting.
- The special role of  $[0, \infty]$  in representing information in an lft.
- The special role of the refinement property leading to the notion of a normal product and an extremely efficient test for interval inclusion.
- The outcome minimisation and information overlap tensor absorption strategies.
- Arbitrary base unbiased exact floating point; provides an especially efficient way to control the flow of information and the size of the integer coefficients. Roughly speaking the size of the integer coefficients increases linearly with the number of emitted digits.
- Arbitrary base biased exact floating point; creates the opportunity to add and subtract two real numbers by the manipulation of symbols and a bounded state variable.
- The flatten algorithm; this brings together many of the above ideas into a complete algorithm.
- Transformation of continued fractions into normal products.

## References

- [1] Algirdas Avizienis. Signed-digit number representations for fast parallel arithmetic. *IRE Transactions on electronic computers*, (10):389–400, 1961.

- [2] H.J. Boehm and R. Cartwright. Exact real arithmetic: Formulating real numbers as functions. In D. Turner, editor, *Research Topics in Functional Programming*, pages 43–64. Addison-Wesley, 1990.
- [3] Claude Brezinski. *History of Continued Fractions and Padé Approximants*, volume 12 of *Springer series in Computational mathematics*. Springer-Verlag, 1991.
- [4] M. H. Escardó. PCF extended with real numbers. *Theoretical Computer Science*, 162(1):79–115, August 1996.
- [5] Pietro Di Gianantonio. Real Number Computability and Domain Theory. In *Proceedings of the 18<sup>th</sup> International Symposium on Mathematical Foundations of Computer Science*, pages 413–422, Gdansk, Poland, September 1993. LNCS 711.
- [6] Pietro Di Gianantonio. A Golden Ratio Notation for the Real Numbers. Technical report, CWI Amsterdam, 1996.
- [7] R. W. Gosper. Continued Fraction Arithmetic. Technical Report HAK-MEM Item 101B, MIT AI MEMO 239, MIT, February 1972. Available from <ftp://ftp.netcom.com/pub/hb/hbaker/hakmem>.
- [8] A. Grzegorzcyk. On the definition of computable real continuous functions. *Fund. Math.*, 44:61–77, 1957.
- [9] IEEE. IEEE Standard 754 for Binary Floating-Point Arithmetic. *SIGPLAN*, 22(2):9–25, 1985.
- [10] Simon L. Peyton Jones. Arbitrary precision arithmetic using continued fractions, 1984. INDRA Note 1530, University College London.
- [11] D. E. Knuth. *The Art of Computer Programming: Seminumerical Algorithms*, volume 2. Addison-Wesley, 1981.
- [12] Peter Kornerup and David W. Matula. An On-line Arithmetic Unit for Bit-Pipelined Rational Arithmetic. *Journal of Parallel and Distributed Computing*, 5(3):310–330, 1988.
- [13] Peter Kornerup and David W. Matula. Exploiting Redundancy in Bit-Pipelined Rational Arithmetic. In *Proceedings of the 9<sup>th</sup> IEEE Symposium on Computer Arithmetic*, pages 119–126, Santa Monica, 1989. IEEE Computer Society Press.
- [14] Peter Kornerup and David W. Matula. An Algorithm for Redundant Binary Bit-Pipelined Rational Arithmetic. *IEEE Transactions on Computers*, C-39(8):1106–1115, 1990.



- [15] Peter Kornerup and David W. Matula. Finite Precision Lexicographic Continued Fraction Number Systems. In *Proceedings of the 7<sup>th</sup> IEEE Symposium on Computer Arithmetic*, pages 207–214, Urbana, 1995. IEEE Computer Society Press.
- [16] V. Menissier-Morain. Arbitrary precision real arithmetic: design and algorithms. submitted to *J. Symbolic Computation*, 1996.
- [17] R.E. Moore. *Interval Analysis*. Prentice-Hall, Englewood Cliffs, 1966.
- [18] Asger Munk Nielsen and Peter Kornerup. MSB-First Digit Serial Arithmetic. *Journal of Universal Computer Science*, 1(7):527–547, July 1995.
- [19] W. Parry. On the  $\beta$ -expansions of real numbers. *Acta Mathematica, Acad. Sci. Hung.*, 11:401–416, 1960.
- [20] P. J. Potts. Computable Real Arithmetic using Linear Fractional Transformations, June 1996. Early draft PhD Thesis, Imperial College, available from <http://www-tfm.doc.ic.ac.uk/~pjp>.
- [21] J. Vuillemin. Exact real computer arithmetic with continued fractions. *IEEE Transactions on computers*, 39(8):1087–1105, August 1990.
- [22] H. S. Wall. *Analytic Theory of Continued Fractions*. Chelsea Publishing Company, 1948.
- [23] Osaaki Watanuki and Milos D. Ercegovic. Error Analysis of Certain Floating-Point On-Line Algorithms. *IEEE Transactions on Computers*, C-32(4):352–358, April 1983.