

# Performing Arithmetic Operations on Round-to-Nearest Representations

Adrien Panhaleux, Peter Kornerup, Jean-Michel Muller

► **To cite this version:**

Adrien Panhaleux, Peter Kornerup, Jean-Michel Muller. Performing Arithmetic Operations on Round-to-Nearest Representations. Pre-print, 24 pages, submitted to IEEE Transactions on Computers. 2009. <ensl-00441933>

**HAL Id: ensl-00441933**

**<https://hal-ens-lyon.archives-ouvertes.fr/ensl-00441933>**

Submitted on 17 Dec 2009

**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.

# Performing Arithmetic Operations on Round-to-Nearest Representations

Peter Kornerup

University of Southern Denmark

Odense, Denmark

*E-mail: kornerup@imada.sdu.dk*

Jean-Michel Muller

CNRS-LIP-Arénaire

Lyon, France

*E-mail: Jean-Michel.Muller@ens-lyon.fr*

Adrien Panhaleux

École Normale Supérieure de Lyon

Lyon, France

*E-mail: Adrien.Panhaleux@ens-lyon.fr*

## Abstract

During any composite computation there is a constant need for rounding intermediate results before they can participate in further processing. Recently a class of number representations denoted RN-Codings were introduced, allowing an un-biased rounding-to-nearest to take place by a simple truncation, with the property that problems with double-roundings are avoided. In this paper we first investigate a particular encoding of the binary representation. This encoding is generalized to any radix and digit set; however radix complement representations for even values of the radix turn out to be particularly feasible. The encoding is essentially an ordinary radix complement representation with an appended round-bit, but still allowing rounding to nearest by truncation and thus avoiding problems with double-roundings. Conversions from radix complement to these round-to-nearest representations can be performed in constant time, whereas conversion the other way in general takes at least logarithmic time. Not only is rounding-to-nearest a constant time operation, but so is also sign inversion, both of which are at best log-time operations on ordinary 2's complement representations. Addition and multiplication on such fixed-point representations are first analyzed and defined in such a way that rounding information can be carried along in a meaningful way, at minimal cost. The analysis is carried through for a compact (canonical) encoding using 2's complement representation, supplied with a round-bit. Based on the fixed-point encoding it is shown possible to define floating point representations, and a sketch of the implementation of an FPU is presented.

## I. INTRODUCTION

In a recent paper [KM05] a class of number representations denoted RN-Codings were introduced, the “RN” standing for “round to nearest”, as these radix- $\beta$ , signed-digit representations have the property that truncation yields rounding to the nearest representable value. They are based on a generalization of the observation that certain radix representations are known to possess this property, e.g., the balanced ternary ( $\beta = 3$ ) system over the digit set  $\{-1, 0, 1\}$ . Another such representation is obtained by performing the original Booth-recoding [Boo51] on

a 2's complement number into the digit set  $\{-1, 0, 1\}$ , where it is well-known that the non-zero digits of the recoded number alternate in sign. To distinguish between situations where we are not concerned with the actual encoding of a value, we shall here use the notation *RN-representation*.

We shall in Section II from [KM05] cite some of the definitions and properties of the general RN-Codings/representations. However, we will in particular explore the binary representation, e.g., as obtained by the Booth recoding; the rounding by truncation property, including the feature that the effect of one rounding followed by another rounding yields the same result, as would be obtained by a single rounding to the same precision as the last.

Section III analyzes conversions between RN-representations and 2's complement representations. Conversion from the latter to the former is performed by the Booth algorithm, yielding a signed-digit/borrow-save representation in a straightforward encoding, which for an  $n$ -digit word requires  $2n$  bits. It is then realized that  $n + 1$  bits are sufficient, providing a simpler alternative encoding consisting of the bits of the truncated 2's complement encoding, with a round-bit appended, termed the *canonical* encoding. Despite being based on a 2's complement encoding, it is observed that sign-inversion (negation) is a constant time operation on this canonical encoding. Conversion the other way, from RN-representation in this encoding into 2's complement representation (essentially adding in the round-bit) is realizable by a parallel prefix structure. Section IV generalizes the canonical representation to other radices and digit sets, showing that for even values of the radix the encodings employing radix-complement representations are particularly feasible.

Section V then analyzes possible implementations of addition and multiplication on fixed-point RN-represented numbers. [BM05] discussed implementations of these basic operations based on the signed-digit representation of RN-coded numbers, whereas we here exploit the canonical encoding, which seems to be more convenient. Since it turns that there are two possible encodings of the result of an arithmetic operation, interpretations of the encodings as intervals may be used to uniquely define sums and products in a consistent way. Section VI sketches how a floating point RN-representation may be defined and the basic arithmetic operations of an FPU may be realized. Then Section VII contains examples on some composite computations where fast and optimal roundings are useful, and may come for free when RN-representation in the canonical encoding is employed. Finally Section VIII concludes the paper.

## II. DEFINITIONS AND BASIC PROPERTIES (CITED FROM [KM05])

*Definition 1 (RN-representations):* Let  $\beta$  be an integer greater than or equal to 2. The digit sequence  $D = d_n d_{n-1} d_{n-2} \dots$  (with  $-\beta + 1 \leq d_i \leq \beta - 1$ ) is an RN-representation in radix  $\beta$  of  $x$  iff

- 1)  $x = \sum_{i=-\infty}^n d_i \beta^i$  (that is  $D$  is a radix- $\beta$  representation of  $x$ );
- 2) for any  $j \leq n$ ,

$$\left| \sum_{i=-\infty}^{j-1} d_i \beta^i \right| \leq \frac{1}{2} \beta^j,$$

that is, if the digit sequence is truncated to the right at any position  $j$ , the remaining sequence is always the number (or one of the two members in case of a tie) of the form  $d_n d_{n-1} d_{n-2} d_{n-3} \dots d_j$  that is closest to  $x$ .

Hence, truncating the RN-representation of a number at any position is equivalent to rounding it to the nearest.

Although it is possible to deal with infinite representations, we shall first restrict our discussions to finite representations. The following observations on such RN-representations for general  $\beta \geq 2$  are then easily found:

*Theorem 2 (Characterizations of finite RN-representations):*

- if  $\beta \geq 3$  is odd, then  $D = d_m d_{m-1} \cdots d_\ell$  is an RN-representation iff

$$\forall i, \frac{-\beta + 1}{2} \leq d_i \leq \frac{\beta - 1}{2};$$

- if  $\beta \geq 2$  is even then  $D = d_m d_{m-1} \cdots d_\ell$  is an RN-representation iff
  - 1) all digits have absolute value less than or equal to  $\frac{\beta}{2}$ ;
  - 2) if  $|d_i| = \frac{\beta}{2}$ , then the first non-zero digit that follows on the right has the opposite sign, that is, the largest  $j < i$  such that  $d_j \neq 0$  satisfies  $d_i \times d_j < 0$ .

Observe that for odd  $\beta$  the system is non-redundant, whereas for  $\beta$  even the system is redundant, in the sense that some non-zero numbers have two representations. In particular note that for radix 2 the digit set is  $\{-1, 0, 1\}$ , known by the names of “signed-digit” or “borrow-save”, but here restricted such that the non-zero digits have alternating signs.

*Theorem 3 (Uniqueness of finite representations):*

- if  $\beta$  is odd, then a finite RN-representation of  $x$  is unique;
- if  $\beta$  is even, then some numbers may have two finite representations. In that case, one has its least significant nonzero digit equal to  $-\frac{\beta}{2}$ , the other one has its least significant nonzero digit equal to  $+\frac{\beta}{2}$ .

*Proof:* If  $\beta$  is odd, the result is an immediate consequence of the fact that the digit set is non-redundant. If  $\beta$  is even, then consider two different RN-representations representing the same value  $x$ , and consider the largest position  $j$  (that is, of weight  $\beta^j$ ) such that these RN-representations differ, when truncated to the right of position  $j$ . Let  $x_a$  and  $x_b$  be the values represented by these digit strings. Obviously,  $x_a - x_b \in \{-\beta^j, 0, \beta^j\}$ . Now  $x_a = x_b$  would contradict the way that  $j$  was chosen. Without loss of generality, then assume  $x_b = x_a + \beta^j$ . This implies  $x = x_a + \beta^j/2 = x_b - \beta^j/2$ , since the maximal absolute value of a digit is  $\beta/2$ . Hence, the remaining digit strings (i.e., the parts that were truncated) are digit strings starting from position  $j - 1$ , representing  $\pm\beta^j/2$ .

The only way of representing  $\beta^j/2$  by an RN-representation starting from position  $j - 1$  is

$$\left(\frac{\beta}{2}\right) 0000 \cdots 0.$$

This is seen as follows: if the digit at position  $j - 1$  of a number is less than or equal to  $\frac{\beta}{2} - 1$ , then that number is less than or equal to

$$\left(\frac{\beta}{2} - 1\right) \beta^{j-1} + \left(\frac{\beta}{2}\right) \sum_{i=\ell}^{j-2} \beta^i < \beta^j/2,$$

since the largest allowed digit is  $\frac{\beta}{2}$ . Also, the digit at position  $j - 1$  of an RN-representation cannot be larger than or equal to  $\frac{\beta}{2} + 1$ . This ends the proof for odd  $\beta$ . ■

If  $\beta$  is even, then a number whose finite representation (by an RN-representation) has its last nonzero digit equal to  $\frac{\beta}{2}$  has an alternative representation ending with  $-\frac{\beta}{2}$  (just assume the last two digits are  $d(\frac{\beta}{2})$ : since the representation is an RN-representation,  $d < \frac{\beta}{2}$ , hence if we replace these two digits by  $(d+1)(-\frac{\beta}{2})$  we still have a valid RN-representation). This has an interesting

consequence: truncating a number which is a tie will round either way, depending on which of the two possible representations the number happens to have. Hence, there is no bias in the rounding.

Note that this rounding rule is different from the “round-to-nearest-even” rule required by the IEEE standard [IEE85]. Both roundings provide a “round-to-nearest” in the case of a tie, but employ different rules when choosing which way to round. Also note that this rounding is also deterministic, the direction of rounding only depends on how the value to be rounded was derived, as the representation of the value is uniquely determined by the sequence of operations leading to the value.

*Example:*

- In radix 7, with digits  $\{-3, -2, -1, 0, 1, 2, 3\}$ , all representations are RN-representations, and no number has more than one representation;
- in radix 10 with digits  $\{-5, \dots, +5\}$ , 15 has two RN-representations: 15 and  $2\bar{5}$ . ■

*Theorem 4 (Uniqueness of infinite representations):* We now consider *infinite* representations, i.e., representations that do not ultimately terminate with an infinite sequence of zeros.

- if  $\beta$  is odd, then some numbers may have two infinite RN-representations. In that case, one is eventually finishing with the infinite digit string

$$\frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \dots$$

and the other one is eventually finishing with the infinite digit string

$$\frac{-\beta+1}{2} \frac{-\beta+1}{2} \frac{-\beta+1}{2} \frac{-\beta+1}{2} \frac{-\beta+1}{2} \frac{-\beta+1}{2} \dots;$$

- if  $\beta$  is even, then two different infinite RN-representations necessarily represent different numbers. As a consequence, a number that is not an integer multiple of an integral (positive or negative) power of  $\beta$  has a unique RN-representation.

*Proof:* If  $\beta$  is odd, the existence immediately comes from

$$1. \frac{-\beta+1}{2} \frac{-\beta+1}{2} \frac{-\beta+1}{2} \frac{-\beta+1}{2} \dots = 0. \frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \dots = \frac{1}{2}$$

Now, if for any  $\beta$  (odd or even) two different RN-representations represent the same number  $x$ , then consider them truncated to the right of some position  $j$ , such that the obtained digit strings differ. The obtained digit strings represent values  $x_a$  and  $x_b$  whose difference is  $\pm\beta^j$  (a larger difference is impossible for obvious reasons).

First, consider the case where  $\beta$  is odd. From the definition of RN-representations, and assuming  $x_a < x_b$ , we have  $x = x_a + \beta^j/2 = x_b - \beta^j/2$ . Since  $\beta$  is odd, the only way of representing  $\beta^j/2$  is with the infinite digit string (that starts from position  $j - 1$ )

$$\frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \frac{\beta-1}{2} \dots$$

the result immediately follows.

Now consider the case where  $\beta$  is even. Let us first show that  $x_a = x_b$  is impossible. From Theorem 3, this would imply that one of the corresponding digit strings would terminate with the digit sequence  $-\frac{\beta}{2}00\dots00$ , and the other one with the digit string  $+\frac{\beta}{2}00\dots00$ . But from Theorem 2, this would imply that the remaining (truncated) terms are positive in the first case, and negative in the second case, which would mean (since  $x_a = x_b$  implies that they are equal) that they would both be zero, which is not compatible with the fact that the representations of  $x$  are assumed infinite. Hence  $x_a \neq x_b$ . Assume  $x_a < x_b$ , which implies  $x_b = x_a + \beta^j$ . We necessarily

have  $x = x_a + \beta^j/2 = x_b - \beta^j/2$ . Although  $\beta^j/2$  has several possible representations in a “general” signed-digit radix- $\beta$  system, the only way of representing it with an RN-representation is to put a digit  $\frac{\beta}{2}$  at position  $j - 1$ , no infinite representation is possible. ■

*Example:*

- In radix 7, with digits  $\{-3, -2, -1, 0, 1, 2, 3\}$ , the number  $3/2$  has two infinite representations, namely  $1.333333333\cdots$  and  $2.\overline{333333333}\cdots$
- in radix 10 with digits  $\{-5, \dots, +5\}$ , the RN-representation of  $\pi$  is unique. ■

An important property of the RN-representation is that it avoids the *double rounding* problem occurring with some rounding methods, e.g., with the standard IEEE round-to-nearest-even. This may happen when the result of first rounding to a position  $j$ , followed by rounding to position  $k$ , does not yield the same result as if directly rounding to position  $k$ . We repeat from [KM05] the following result:

*Observation 5 (Double rounding):*

Let  $\text{rn}_i(x)$  be the function that rounds the value of  $x$  to nearest at position  $i$  by truncation. Then for  $k > j$ , if  $x$  is represented in the RN-representation, then

$$\text{rn}_k(x) = \text{rn}_k(\text{rn}_j(x))$$

### III. CONVERTING TO AND FROM BINARY RN-REPRESENTATION

#### A. Conversion from 2's Complement to RN-Representation

Consider an input value  $x = -b_m 2^m + \sum_{i=\ell}^{m-1} b_i 2^i$  in 2's complement representation:

$$x \sim b_m b_{m-1} \cdots b_{\ell+1} b_\ell$$

with  $b_i \in \{0, 1\}$  and  $m > \ell$ . Then the digit string

$$\delta_m \delta_{m-1} \cdots \delta_{\ell+1} \delta_\ell \quad \text{with} \quad \delta_i \in \{-1, 0, 1\}$$

defined (by the Booth recoding [Boo51]) for  $i = \ell, \dots, m$  as

$$\delta_i = b_{i-1} - b_i \quad (\text{with } b_{\ell-1} = 0 \text{ by convention}) \quad (1)$$

is an RN-representation of  $x$  with  $\delta_i \in \{-1, 0, 1\}$ . That it represents the same value follows trivially by observing that the converted string represents the value  $2x - x$ . The alternation of the signs of non-zero digits is easily seen by considering how strings of the form  $011\cdots 10$  and  $100\cdots 01$  are converted.

Thus the conversion can be performed in constant time. Actually, the digits of the 2's complement representation directly provides for an encoding of the converted digits as a tuple:  $\delta_i \sim (b_{i-1}, b_i)$  for  $i = \ell, \dots, m$  where

$$\begin{aligned} -1 &\sim (0, 1) \\ 0 &\sim (0, 0) \text{ or } (1, 1) \\ 1 &\sim (1, 0), \end{aligned} \quad (2)$$

where the value of the digit is the difference between the first and the second component.

*Example:* Let  $x = 110100110010$  be a sign-extended 2's complement number and write the digits of  $2x$  above the digits of  $x$ :

$2x$	1	0	1	0	0	1	1	0	0	1	0	0
$x$	1	1	0	1	0	0	1	1	0	0	1	0
$x$ in RN-representation		$\bar{1}$	1	$\bar{1}$	0	1	0	$\bar{1}$	0	1	$\bar{1}$	0

where it is seen that in any column the two upper-most bits provide the encoding defined above of the signed-digit below in the column. Since the digit in position  $m+1$  will always be 0, there is no need to include the most significant position otherwise found in the two top rows. ■

If  $x$  is non-zero and  $b_k$  is the least significant non-zero bit of the 2' complement representation of  $x$ , then  $\delta_k = -1$ , confirmed in the example, hence the last non-zero digit is always  $\bar{1}$  and thus unique. However, if an RN-represented number is truncated for rounding somewhere, the resulting representation may have its last non-zero digit of value 1.

As mentioned in Theorem 3 there are exactly two finite binary RN-representations of any non-zero binary number of the form  $a2^k$  for integral  $a$  and  $k$ , but requiring a specific sign of the last non-zero digit makes the representation unique. On the other hand without this requirement, rounding by truncation makes the rounding unbiased in the tie-situation, by randomly rounding up or down, depending on the sign of the last non-zero digit in the remaining digit string.

*Example:* Rounding the value of  $x$  in Example 1 by truncating off the two least significant digits we obtain

$\text{rn}_2(2x)$	1	0	1	0	0	1	1	0	0	<b>1</b>
$\text{rn}_2(x)$	1	1	0	1	0	0	1	1	0	0
$\text{rn}_2(x)$ in RN-representation		$\bar{1}$	1	$\bar{1}$	0	1	0	$\bar{1}$	0	1

where it is noted that the bit of value 1 in the upper rightmost corner (in boldface) acts as a round bit, assuring a round-up in cases there is a tie-situation as here. ■

The example shows that there is another very compact encoding of RN-represented numbers derived directly from the 2's complement representation, noting in the example that the upper row need not be part of the encoding, except for the round-bit. We will denote it the *canonical encoding*, and note that it is a kind of "carry-save" in the sense that it contains a bit not yet added in. The same idea have previously been pursued in [NMLE00] in a floating-point setting, denoted "packet-forwarding".

*Definition 6 (Binary canonical RN-encoding):*

Let the number  $x$  be given in 2's complement representation as the bit string  $b_m \cdots b_{\ell+1} b_\ell$ , such that  $x = -b_m 2^m + \sum_{i=\ell}^{m-1} b_i 2^i$ . Then the *binary canonical encoding* of the RN-representation of  $x$  is defined as the pair

$$x \sim (b_m b_{m-1} \cdots b_{\ell+1} b_\ell, r) \quad \text{where the round-bit is } r = 0$$

and after truncation at position  $k$ , for  $m \geq k > \ell$

$$\text{rn}_k(x) \sim (b_m b_{m-1} \cdots b_{k+1} b_k, r) \quad \text{with round-bit } r = b_{k-1}.$$

If  $(x, r_x)$  is the binary canonical (2's complement) RN-representation of  $X$ , then  $X = x + r_x u$ , from which it follows that

$$-X = -x - r_x u = \bar{x} + u - r_x u = \bar{x} + (1 - r_x)u = \bar{x} + \bar{r}_x u.$$

*Observation 7:* If  $(x, r_x)$  is the canonical RN-representation of a value  $X$ , then  $(\bar{x}, \bar{r}_x)$  is the canonical RN-representation of  $-X$ , where  $\bar{x}$  is the 1's complement of  $x$ . Hence negation of a canonically encoded value is a constant time operation.

The signed-digit interpretation is available from the canonical encoding by pairing bits,  $(b_{i-1}, b_i)$  using the encoding (2) for  $i > k$  and  $(r, b_k)$ , when truncated at position  $k$ .

There are other equally compact encodings of RN-represented numbers, e.g., one could encode the signed-digit string simply by the string of bits obtained as the absolute values of the digits, together with say the sign of the most (or least) non-zero digit. Due to the alternating signs of the non-zero digits, this is sufficient to reconstruct the actual digit values. However, this encoding does not seem very convenient for arithmetic processing, as the correct signs will then have to be distributed over the bit string.

### B. Conversion from Signed-Digit RN-Representation to 2's Complement

The example of converting  $0000000\bar{1}$  into its 2's complement equivalent  $11111111$  shows that it is not possible to perform this conversion in constant time, information may have to travel an arbitrary distance to the left. Hence a conversion may in general take at least logarithmic time. Since the RN-representation is a special case of the (redundant) signed-digit representation, this conversion is fundamentally equivalent to an addition.

If an RN-represented number is in canonical encoding, conversion into ordinary 2's complement representation may require a non-zero round-bit to be added in, it simply consists in an incrementation, for which very efficient methods exist based on parallel prefix trees with AND-gates as nodes.

## IV. CANONICAL REPRESENTATION, THE GENERAL CASE

The binary canonical representation of RN-representation is specified by  $x = (a, r_a)$ , which is a pair of a number and a bit. We could decide to represent the value of  $a$  of that pair in something else than binary, say using a higher radix  $\beta$  and/or a digit set different from the set  $\{0, \dots, \beta - 1\}$ .

*Definition 8 (Canonical encoding: general case):* Let  $b$  be a number in radix  $\beta$  using the digit set  $D$ , such that  $b = \sum_{i=\ell}^{m-1} b_i \beta^i$  with  $b_i \in D$ , and the rounding bit  $r_b \in \{0, 1\}$ . The pair  $(b, r_b)$  then represents the value  $b + ur_b$ , where  $u$  is the unit in the last place ( $u = \beta^\ell$ ).

The definition is very general, as the representation doesn't necessarily allow rounding by truncation. We must redefine the rounding operation so that we avoid problems with double-roundings, basically by trying to convert the encoding into an RN-representation satisfying Definition 1.

### A. Even Radix

The definition seems to make particular sense when a (non-negative) number is represented in an even radix with the regular digit-set  $\{0, \dots, \beta - 1\}$ . In that representation the link between the canonical encoding and RN-representation is trivial enough, so that rounding-to-nearest can be done by truncation of the value  $b$  in the pair  $(b, r_b)$ .

Consider an input value in radix- $\beta$  with  $0 \leq d_i \leq \beta - 1$

$$(x, r_x) = (d_m d_{m-1} d_{m-2} \cdots d_\ell, r_x),$$

and define variables  $c_k$  as

$$c_{k+1} = \begin{cases} 1 & \text{if } d_k \geq \frac{\beta}{2} \\ 0 & \text{if } d_k < \frac{\beta}{2} \end{cases} \quad (3)$$

With  $c_\ell = r_x$  the digits  $\delta_k$  of an RN-representation can be obtained using

$$\delta_k = d_k + c_k - \beta c_{k+1}.$$



The conversion for an even radix and digit set  $\{0, \dots, \beta - 1\}$  into RN-representation gives us a way to easily perform rounding-to-nearest by truncation of  $(x, r_x)$  in the canonical encoding. For  $k > \ell$ :

$$\text{rn}_k(x, r_x) \sim (d_n d_{n-1} \dots d_{k+1} d_k, r) \quad \text{with round-bit } r = \begin{cases} 1 & \text{if } d_{k-1} \geq \beta/2 \\ 0 & \text{if } d_{k-1} < \beta/2 \end{cases}$$

so in RN-representation the value can also be expressed by the digit string  $\delta_n \delta_{n-1} \dots \delta_k \in \{-\frac{\beta}{2}, \dots, \frac{\beta}{2}\}$ .

*Example:* With radix  $\beta = 10$  and the regular digit-set  $\{0, \dots, 9\}$ , for the value 9.25451 represented by  $(9.25450, 1)$ , we can truncate using the previous algorithm:

$$\text{rn}_{-3}(9.25450, 1) = (9.254, 1), \text{ the rounding-bit being } 1 \text{ because } d_{-4} = 5 \implies c_{-3} = 1.$$

We only need to generate one carry ( $c_{-3}$ ) to obtain the rounding bit.

To confirm that the rounding is correct, we may represent the value in RN-representation, by generating all the carries:

$c_k$	1	0	1	0	1	0	<b>1</b>
$d_k$	0	9	2	5	4	5	0
RN-representation	1	$\bar{1}$	3	$\bar{5}$	5	$\bar{5}$	1

With this RN-representation, truncating at position  $-3$  gives:

$c_k$	1	0	1	0	<b>1</b>
$d_k$	0	9	2	5	4
RN-representation	1	$\bar{1}$	3	$\bar{5}$	5

corresponding to the previous canonical encoding  $(9.254, 1)$ . Note that to represent the given positive value,  $d_1$  was set to zero. Had the value been a (negative) 10's complement represented number, then  $d_1$  should by sign extension have been set to 9. ■

### B. Other Representations

If we use other representations of  $b$  (say binary borrow-save, odd radices,...), the rounding may take time  $O(\log(n))$ :

*Example:* [Borrow-save:] When trying to round  $x$  in a general borrow-save representation to the nearest integer, we have for any round bit  $r_x \in \{0, 1\}$ :

borrow-save	rounded
$(1\bar{1}\bar{1}0.00 \dots 0\bar{1}, r_x)$	$(1\bar{1}\bar{1}\bar{1}, 1)$
$(1\bar{1}\bar{1}0.00 \dots 00, r_x)$	$(1\bar{1}\bar{1}0, 0)$

hence we may have to look arbitrarily far to the right when rounding the values. ■

However, borrow-save could be interesting, since addition then can be performed in  $O(1)$ , instead of  $O(\log(n))$  for binary canonical encoding using the regular digit-set  $\{0, 1\}$ . It is important to recall that borrow-save isn't an RN-representation, even though it uses the same digits. To have an RN-representation, the non-zero digits must alternate in signs, and translating an arbitrary number from borrow-save to RN-representation might take a  $O(\log(n))$  time.

*Example:* [Odd radices:] When trying to round to the nearest integer, we have similarly:

radix 3	rounded
$(10.11 \dots 12, r_x)$	$(10, 1)$
$(10.11 \dots 11, r_x)$	$(10, 0)$

which means we may have to look arbitrarily far to the right when rounding the values. It is due to the fact that the mid-point between two representable numbers needs to be represented with an infinite number of digits. In general, odd radices do not seem convenient to use. ■

## V. PERFORMING ARITHMETIC OPERATIONS ON CANONICALLY REPRESENTED VALUES

The fundamental idea of the canonical radix-2 RN-representation is that it is a binary representation of some value using the digit set  $\{-1, 0, 1\}$ , but such that non-zero digits alternate in sign. We then introduced an encoding of such numbers, employing 2's complement representation, in the form  $(a, r_a)$  representing the value

$$(2a + r_a u) - a = a + r_a u,$$

where  $u$  is the weight of the least significant position of  $a$ . Note that there is then no difference between  $(a, 1)$  and  $(a + u, 0)$ , both being RN-representations of the same value:

$$\forall a, \mathcal{V}(a, 1) = \mathcal{V}(a + u, 0),$$

where we use the notation  $\mathcal{V}(x, r_x)$  to denote the value of an RN-represented number.

### A. An Interval Interpretation

Considered as intervals as described below, the two representations  $(a, 1)$  and  $(a+u, 0)$  describe different intervals. Since different representations of the same number can give different rounding results when truncated, it is then important to choose carefully the representation of the result when performing arithmetic operations like addition and multiplication. Hence when defining the result it is essential to choose the encoding of it to reflect the domains of the operands.

Consider a value  $A$  to be rounded at some position of weight  $u$  where the round bit is 1, shown in boldface:

$$\begin{array}{r|l} \dots 0 \mathbf{1} 1 \dots \mathbf{1} & \text{x} \dots \\ \dots 0 1 \dots 1 & 1 \text{x} \dots \\ \hline \underbrace{\dots 1 0 \dots 0}_{a+u} & \underbrace{\dots}_{-\frac{u}{2} \leq t \leq 0} \end{array} \qquad \begin{array}{r|l} \dots 0 \mathbf{1} & \text{x} \dots \\ \dots 0 1 & \text{x} \dots \\ \hline \underbrace{\dots 1}_{a+u} & \underbrace{\dots}_{-\frac{u}{2} \leq t \leq 0} \end{array}$$

and similarly when the round bit is 0:

$$\begin{array}{r|l} \dots 1 0 0 \dots \mathbf{0} & \text{x} \dots \\ \dots 1 0 \dots 0 & 0 \text{x} \dots \\ \hline \underbrace{\dots 1 0 \dots 0}_a & \underbrace{\dots}_{0 \leq t \leq \frac{u}{2}} \end{array} \qquad \begin{array}{r|l} \dots 1 \mathbf{0} & \text{x} \dots \\ \dots 1 0 & \text{x} \dots \\ \hline \underbrace{\dots 1}_a & \underbrace{\dots}_{0 \leq t \leq \frac{u}{2}} \end{array}$$

expressing bounds on the tail  $t$  thrown away during rounding by truncation. Observe that the right-hand ends of the intervals are closed, corresponding to a possibly infinite sequence of units

having been thrown away. We find that the value  $A$  before rounding into  $\mathcal{V}(a, r_a)$  must belong to the interval:

$$A \in \left\{ \begin{array}{ll} \left[ a ; a + \frac{u}{2} \right] & \text{for } r_a = 0 \\ \left[ a + \frac{u}{2} ; a + u \right] & \text{for } r_a = 1 \end{array} \right\} = \left[ a + r_a \frac{u}{2} ; a + (1 + r_a) \frac{u}{2} \right] = \mathcal{I}(a, r_a). \quad (4)$$

In the following we shall use  $\mathcal{I}(a, r_a)$  to denote the interval, the idea being to remember where the real number was before rounding.

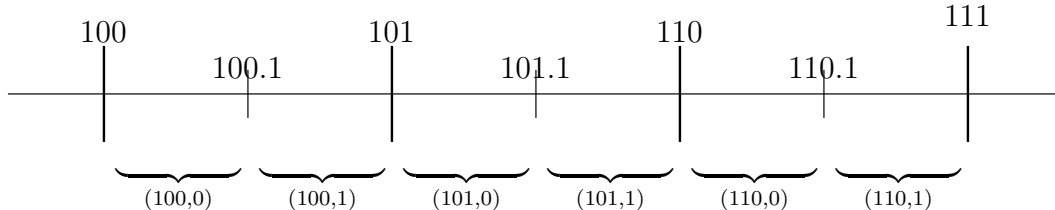


Fig. 1. Example of interpreting RN representations as intervals with  $u = 1$

We may interpret the representations of an encoding as an interval of length  $u/2$ , as in Fig. 1. In the figure, any number between 101 and 101.1 (for example 101.01), when rounded to the nearest integer, will give the RN representation (101, 0). So we may say that (101, 0) represents the interval  $[101; 101.1]$  and in particular

$$\begin{aligned} \mathcal{I}(a, 1) &= \left[ a + \frac{u}{2} ; a + u \right], \\ \mathcal{I}(a + u, 0) &= \left[ a + u ; a + \frac{3u}{2} \right]. \end{aligned}$$

Hence even though the two encodings represent the same value ( $a + u$ ), when interpreting them as intervals according to what could have been thrown away, the intervals are essentially disjoint, except for sharing a single point. In general we may express the interval interpretation as pictured in Fig. 2

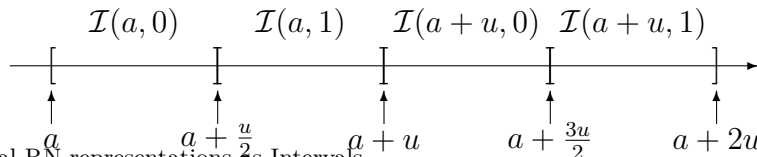


Fig. 2. Binary Canonical RN-representations as Intervals

We do not intend to define an interval arithmetic, but only require that the interval representation of the result of an arithmetic operation  $\odot$  satisfies<sup>1</sup>

$$\mathcal{I}(A \odot B) \subseteq \mathcal{I}(A) \odot \mathcal{I}(B) = \{a \odot b | a \in A, b \in B\}.$$

To simplify the discussion, we will in this section consider only fixed-point representations for some fixed value of  $u$ . We will not discuss overflow problems, as we assume that we have enough bits to represent the result in canonically encoded representation.

<sup>1</sup>Note that this is the reverse inclusion of that required for ordinary interval arithmetic

### B. Addition of RN-represented Values

Employing the value interpretation we have for addition:

$$\begin{array}{r} \mathcal{V}(a, r_a) = a + r_a u \\ + \mathcal{V}(b, r_b) = b + r_b u \\ \hline \mathcal{V}(a, r_a) + \mathcal{V}(b, r_b) = a + b + (r_a + r_b)u \end{array}$$

The resulting value has two possible representations, depending on the rounding bit of the result. To determine what the rounding bit of the result should be, we may look at the interval interpretations of the two possible representations of the result, depending on the rounding bits of the operands.

	$\mathcal{I}(\mathcal{V}(a, r_a) + \mathcal{V}(b, r_b))$	$\mathcal{I}(a, r_a) + \mathcal{I}(b, r_b)$
$\mathbf{r_a = r_b = 0}$	$\mathcal{I}(a + b - u, 1) = \left[ a + b - \frac{u}{2}; a + b \right]$ $\mathcal{I}(a + b, 0) = \left[ a + b; a + b + \frac{u}{2} \right]$	$\not\subseteq [a + b; a + b + u]$ $\subseteq [a + b; a + b + u]$
$\mathbf{r_a \oplus r_b = 1}$	$\mathcal{I}(a + b, 1) = \left[ a + b + \frac{u}{2}; a + b + u \right]$ $\mathcal{I}(a + b + u, 0) = \left[ a + b + u; a + b + \frac{3u}{2} \right]$	$\subseteq \left[ a + b + \frac{u}{2}; a + b + \frac{3u}{2} \right]$
$\mathbf{r_a = r_b = 1}$	$\mathcal{I}(a + b + u, 1) = \left[ a + b + \frac{3u}{2}; a + b + 2u \right]$ $\mathcal{I}(a + b + 2u, 0) = \left[ a + b + 2u; a + b + \frac{5u}{2} \right]$	$\subseteq [a + b + u; a + b + 2u]$ $\not\subseteq [a + b + u; a + b + 2u]$

Since we want  $\mathcal{I}(\mathcal{V}(a, r_a) + \mathcal{V}(b, r_b)) \subseteq \mathcal{I}(a, r_a) + \mathcal{I}(b, r_b)$ , and  $(a, r_a) + (0, 0) = (a, r_a)$ , and in order to keep the addition symmetric, we define the addition of RN encoded numbers as follows.

*Definition 9 (Addition of canonically encoded numbers):* If  $u$  is the unit in the last place of the operands, let:

$$(a, r_a) + (b, r_b) = ((a + b + (r_a \wedge r_b)u), r_a \vee r_b)$$

Noticing that  $-(x, r_x) = (\bar{x}, \bar{r}_x)$ , we observe that using this definition,  $(x, r_x) - (x, r_x) = (-u, 1)$ , with  $\mathcal{V}(-u, 1) = 0$ . In fact, it is possible to define another addition on RN-encoded numbers as  $(a, r_a) +_2 (b, r_b) = ((a + b + (r_a \vee r_b)u), r_a \wedge r_b)$ . Using that alternative definition,  $(x, r_x) -_2 (x, r_x) = (0, 0)$ , but the neutral element for addition is then  $(-u, 1)$ , i.e.  $(x, r_x) +_2 (-u, 1) = (x, r_x)$ .

*Example:* Let us take two examples adding two numbers that were previously rounded to the nearest integer.

	Addition not rounded	Addition on rounded canonical representations
$a_1$	01011.1110	(01011, 1)
$b_1$	01001.1101	(01001, 1)
$a_1 + b_1$	010101.1011	(010101, 1)
$a_2$	01011.1010	(01011, 1)
$b_2$	01001.1001	(01001, 1)
$a_2 + b_2$	010101.0011	(010101, 1)

Using the definition above,  $\text{rn}_0(a_1 + b_1) = \text{rn}_0(a_1) + \text{rn}_0(b_1)$  holds in the first case. Obviously, since some information may be lost during rounding, there are cases like in the second example where  $\text{rn}_0(a_2 + b_2) \neq \text{rn}_0(a_2) + \text{rn}_0(b_2)$ . Also note that due to that information loss,  $a_2 + b_2$  is not in  $\mathcal{I}((a_2, r_{a_2}) + (b_2, r_{b_2}))$   $\blacksquare$

When interpreted as an interval,  $\mathcal{I}(x, r_x) = \left[ x + r_x \frac{u}{2}; x + (1 + r_x) \frac{u}{2} \right]$  then its “mirror image” interval of negated values is  $-\mathcal{I}(x, r_x) = \left[ \bar{x} + \bar{r}_x \frac{u}{2}; \bar{x} + (1 + \bar{r}_x) \frac{u}{2} \right] = \mathcal{I}(\bar{x}, \bar{r}_x)$ . Thus consider for subtraction

$$\begin{aligned}
\mathcal{I}((a, r_a) - (b, r_b)) &= \mathcal{I}((a, r_a) + (\bar{b}, \bar{r}_b)) \\
&= \mathcal{I}(a + \bar{b} + (r_a \vee \bar{r}_b)u, r_a \wedge \bar{r}_b) \\
&= \left[ a + \bar{b} + u(r_a \vee \bar{r}_b) + \frac{u}{2}(r_a \wedge \bar{r}_b); a + \bar{b} + u(r_a \vee \bar{r}_b) + \frac{u}{2}(1 + r_a \wedge \bar{r}_b) \right] \\
&= a + \bar{b} + u(r_a \vee \bar{r}_b) + \frac{u}{2}(r_a \wedge \bar{r}_b) + \left[ 0; \frac{u}{2} \right] \\
&= a - b + u(r_a - r_b) - \frac{u}{2}(r_a \wedge \bar{r}_b) + \left[ 0; \frac{u}{2} \right] \\
&= a - b + \frac{u}{2}(r_a - r_b) - \frac{u}{2}(\bar{r}_a \wedge r_b) + \left[ 0; \frac{u}{2} \right].
\end{aligned}$$

On the other hand,

$$\begin{aligned}
\mathcal{I}(a, r_a) - \mathcal{I}(b, r_b) &= \left[ a + \frac{u}{2}r_a; a + \frac{u}{2}(1 + r_a) \right] - \left[ b + \frac{u}{2}r_b; b + \frac{u}{2}(1 + r_b) \right] \\
&= a - b + \frac{u}{2}(r_a - r_b) + \left[ -\frac{u}{2}; \frac{u}{2} \right],
\end{aligned}$$

hence for all points  $x \in \mathcal{I}((a, r_a) - (b, r_b))$ ,  $x$  is in  $\mathcal{I}(a, r_a) - \mathcal{I}(b, r_b)$ .

Hence subtraction of  $(x, r_x)$  can be realized by addition of the bitwise inverted tuple  $(\bar{x}, \bar{r}_x)$ .

### C. Multiplying Canonically Encoded RN-represented values

By definition we have for the value of the product

$$\begin{array}{rcl}
\mathcal{V}(a, r_a) &= & a + r_a u \\
\mathcal{V}(b, r_b) &= & b + r_b u \\
\hline
\mathcal{V}(a, r_a)\mathcal{V}(b, r_b) &= & ab + (ar_b + br_a)u + r_a r_b u^2,
\end{array}$$

noting that the unit of the result is  $u^2$ , assuming that  $u \leq 1$ . Considering the operands as intervals we find using (4):

$$\begin{aligned}
\mathcal{I}(a, r_a) \times \mathcal{I}(b, r_b) &= \left[ a + r_a \frac{u}{2}; a + \frac{u}{2}(1 + r_a) \right] \times \left[ b + \frac{u}{2}r_b; b + \frac{u}{2}(1 + r_b) \right] \\
&= a'b' + \begin{cases} \left[ 0; a' + b' + \frac{u}{2} \right] \times \frac{u}{2} & \text{for } a > 0, b > 0 \\ \left[ a'; b' \right] \times \frac{u}{2} & \text{for } a < 0, b > 0 \\ \left[ b'; a' \right] \times \frac{u}{2} & \text{for } a > 0, b < 0 \\ \left[ a' + b' + \frac{u}{2}; 0 \right] \times \frac{u}{2} & \text{for } a < 0, b < 0 \end{cases} \quad (5)
\end{aligned}$$

where  $a' = a + r_a \frac{u}{2}$  and  $b' = b + r_b \frac{u}{2}$ , and it is assumed that  $a$  and  $a'$  share the same sign and similarly for  $b$  and  $b'$  (assumed satisfied if  $a \neq 0$  and  $b \neq 0$ ). But for  $a < 0$  and  $b < 0$ , with  $r_a = r_b = 0$  we would expect the result  $(ab, 0)$  as interval  $\mathcal{I}(ab, 0) = ab + \left[ 0; \frac{u}{2} \right]$ , to be in the appropriate interval defined by (5), which is NOT the case!

However, since negation of canonical (2's complement) RN-represented values can be obtained by constant-time bit inversion, multiplication of such operands can be realized by multiplication of the absolute values of the operands, the result being supplied with the correct sign by a conditional inversion.

Thus employing bit-wise inversions, multiplication in 2's complement RN-representations becomes equivalent to sign-magnitude multiplication, hence assuming that both operands are non-negative, the "interval product" is

$$\begin{aligned}
\mathcal{I}(a, r_a) \times \mathcal{I}(b, r_b) &= \left[ a + r_a \frac{u}{2} ; a + \frac{u}{2}(1 + r_a) \right] \times \left[ b + r_b \frac{u}{2} ; b + \frac{u}{2}(1 + r_b) \right] \\
&= \left[ (a + r_a \frac{u}{2})(b + r_b \frac{u}{2}) ; (a + \frac{u}{2}(1 + r_a))(b + \frac{u}{2}(1 + r_b)) \right] \\
&= (a + r_a \frac{u}{2})(b + r_b \frac{u}{2}) + \left[ 0 ; (a + r_a \frac{u}{2}) + (b + r_b \frac{u}{2}) + \frac{u^2}{4} \right] \\
&= ab + \left[ ar_b + br_a + \frac{r_a r_b}{2} u \right] \frac{u}{2} + \left[ 0 ; a + b + (r_a + r_b + 1) \frac{u}{2} \right] \frac{u}{2} \quad (6) \\
&= \begin{cases} ab + \left[ 0 ; a + b + \frac{u}{2} \right] \frac{u}{2} & \text{for } \mathbf{r}_a = \mathbf{r}_b = \mathbf{0} \\ ab + a \frac{u}{2} + \left[ 0 ; a + b + u \right] \frac{u}{2} & \text{for } \mathbf{r}_a = \mathbf{0}, \mathbf{r}_b = \mathbf{1} \\ ab + b \frac{u}{2} + \left[ 0 ; a + b + u \right] \frac{u}{2} & \text{for } \mathbf{r}_a = \mathbf{1}, \mathbf{r}_b = \mathbf{0} \\ ab + (a + b + \frac{u}{2}) \frac{u}{2} + \left[ 0 ; a + b + \frac{3}{2}u \right] \frac{u}{2} & \text{for } \mathbf{r}_a = \mathbf{r}_b = \mathbf{1} \end{cases}
\end{aligned}$$

It then follows that  $\mathcal{I}((ab + (ar_b + br_a)u, r_a r_b) \subseteq \mathcal{I}((a, r_a) \times (b, r_b))$  with unit  $u^2$ , since the lefthand RN-representation corresponds to the interval

$$\left[ (ab + (ar_b + br_a)u) + (r_a r_b) \frac{u^2}{2} ; (ab + (ar_b + br_a)u) + (1 + r_a r_b) \frac{u^2}{2} \right]$$

and its lower endpoint is greater than or equal to the lower endpoint from (6):

$$ab + (ar_b + br_a)u + (r_a r_b) \frac{u^2}{2} \geq ab + (ar_b + br_a + \frac{r_a r_b}{2} u) \frac{u}{2}$$

together with the upper endpoint being smaller than or equal to that from (6):

$$ab + (ar_b + br_a)u + (1 + r_a r_b) \frac{u^2}{2} \leq ab + (ar_b + br_a + \frac{r_a r_b}{2} u) \frac{u}{2} + (a + b + (r_a + r_b + 1) \frac{u}{2}) \frac{u}{2}$$

both satisfied for  $a \geq u$ ,  $b \geq u$  (i.e., non-zero) and all permissible values of  $r_a, r_b$ .

*Definition 10 (Multiplication of canonically encoded numbers):* If  $u$  is the unit in the last place, with  $u \leq 1$ , we define for non-negative operands:

$$(a, r_a) \times (b, r_b) = (ab + u(ar_b + br_a), r_a r_b),$$

and for general operands by appropriate sign inversions of the operands and result. If  $u < 1$  the unit is  $u^2 < u$  and the result may often have to be rounded to unit  $u$ , which can be done by truncation.

For an implementation some modifications to an unsigned multiplier will handle the  $r_a$  and  $r_b$  round bits, we just have to calculate the double length product with two additional rows consolidated into the partial product array. However, we shall not here go into the details of the consolidation.

The multiplier  $(b + r_b u)$  may also be recoded into radix 2 (actually it is already so when interpreted as a signed-digit number) or into radix 4, and a term (bit)  $d_i r_a$  may be added in the row below the row for the partial product  $d_i a$ , where  $d_i$  is the recoded  $i$ 'th digit of the multiplier. Hence only at the very bottom of the array of partial products will there be a need for adding in an extra bit as a new row. The double length product can be returned as  $(p, r_p)$ , noticing that the unit now is  $u' = u^2$ , but the result may have to be rounded, which by simple truncation will define the rounded product as some  $(p', r'_p)$ .

*Example:* As an example with  $u = 1$ :

	Not rounded	Canonical Representation
$a$	01011.1110	(01011, 1)
$b$	01001.1101	(01001, 1)
$a \times b$	01110100.10000110	(01110111, 1)

The multiplication in canonical representation was done according to the definition:

$$ab + (ar_b + br_a) = 01100011 + (01011 + 01001) = 01100011 + 010100 = 01110111,$$

where we note that (01110111, 1) corresponds to the interval:

$$[01110111.1 ; 01111000.0]$$

clearly a subset of the interval

$$[01011.1 \times 01001.1 ; 01100 \times 01010] = [01101101.01 ; 01111000.00].$$

It is obvious that rounding results in larger errors when performing multiplication. ■

Similarly, some other arithmetic operations like squaring, square root or even the evaluation of “well behaved” transcendental functions may be defined and implemented, just considering canonical RN-represented operands as 2’s complement values with a “carry-in” not yet absorbed, possibly using interval interpretation to define the resulting round bit.

## VI. FLOATING POINT REPRESENTATIONS

For an implementation of a floating point arithmetic unit (FPU) it is necessary to define a binary encoding, which we assume is based on the canonical 2’s complement for the encoding of the significand part (say  $s$  encoded in  $p$  bits, 2’s complement), supplied with the round bit (say  $r_s$ ) and an exponent (say  $e$  in some biased binary encoding). It then seems natural to pack the three parts into a computer word (32, 64 or 128 bits) in the following order:



with the round bit in immediate continuation of the significand part, thus simplifying the rounding by truncation. As usual we will require the value being represented is in normalized form, say such that the radix point is between the first and second bit of the significand field. If the first bit is zero, the significand field then represents a fixed point value in the interval  $\frac{1}{2} \leq s < 1$ , if it is one then  $-1 \leq s < -\frac{1}{2}$ .

We shall now sketch how the fundamental operations may be implemented on such floating point RN-representations, not going into details on overflow, underflow and exceptional values.

### A. Multiplication

Since the exponents are handled separately, forming the product of the significands is precisely as described previously for fixed point representations: sign-inverting negative operands, forming the double-length product, normalizing and rounding it, and possibly negating the result, supplying it with the proper sign.

Normalizing here may require a left shift, which is straight forward on the (positive) product before rounding by truncation.

### B. Addition, the "far case"

Addition is traditionally now handled in an FPU as two cases, the "near case" dealing with effective subtraction of operands whose exponents differ by no more than one. Here a significant cancellation of leading digits may occur, and thus a variable amount of left-normalization is required. We will deal with this case below

The remaining cases dealt with are the "far case", where the result at most requires normalization by a single right shift. Otherwise addition/subtraction takes place as for the similar operation in IEEE, sign magnitude representation. There is no need in general to form the exact sum/difference when there is a great difference in exponents.

### C. Subtraction, the "near case"

After cancellation of leading digits there is a need to left normalize, so our problem here is to consider what to shift in from the right. Thinking of the value as represented in signed digit, binary value, obviously zeroes have to be shifted in.

In our encoding, say for a positive result  $(d, r_d)$  we may have a 2's complement bit pattern:

$$d \sim 0 0 \cdots 0 1 b_{k+2} \cdots b_{p-1} \text{ and round bit } r_d$$

to be left normalized. Here the least significant digit is encoded as  $\left\{ \begin{matrix} r_d \\ b_{p-1} \end{matrix} \right\}$ .

It is then found that shifting in bits of value  $r_d$  will precisely achieve the effect of shifting in zeroes in the signed-digit interpretation:

$$2^k d \sim 0 1 b_{k+2} \cdots b_{p-1} r_d \cdots r_d \text{ with round bit } r_d,$$

as seen from  $2 \times (x, r_x) = (x, r_x) + (x, r_x) = (2x + r_x u, r_x)$ .

### D. Division

As for multiplication we assume that negative operands have been sign-inverted, and that exponents are treated separately.

Employing our interval interpretation, we must require the result of division of  $(x, r_x)$  by  $(y, r_y)$  to be in the interval:

$$\left[ \frac{x + r_x \frac{u}{2}}{y + (1 + r_y) \frac{u}{2}}; \frac{x + (1 + r_x) \frac{u}{2}}{y + r_y \frac{u}{2}} \right].$$

After some algebraic manipulations it is found that the exact rational

$$q = \frac{x + r_x u}{y + r_y u}$$

belongs to that interval. Hence any standard division algorithm may be used to develop an approximation to the value of  $q$  to  $(p+1)$ -bit precision, i.e., including the usual round bit where the sign of the remainder may be used to determine if the exact result is just below or above the found approximation.



### E. Discussion of Floating Point RN-Representations

As seen above it is feasible to define a floating point representation where the significand is encoded in the binary canonical 2's complement encoding, together with the round-bit appended at the end of the encoding of the significand. An FPU implementation of the basic arithmetic operations is feasible at about the same complexity as those based on the IEEE-754 standard for binary floating point, with a possible slight overhead in multiplication due to extra terms to be added, but since the round-to-nearest functionality is achieved at much less hardware complexity, the arithmetic operations will generally be faster, by avoiding the usual log-time rounding. The other (directed) roundings can also be realized at minimal cost.

Benefits are obtained through faster rounding and sign inversion (both constant time), but the domain of representable values is not symmetric, as is the case when encoding the significand in sign-magnitude representations.

## VII. APPLICATIONS IN SIGNAL PROCESSING

Let us here think of the use of the fixed-point RN representation in high-speed digital signal processing applications, although there are similar benefits in floating point.

Two particular applications needing frequent roundings come to mind: calculation of inner products for filtering, and polynomial evaluations for approximation of standard functions. For the latter application, a very efficient way of evaluating a polynomial is to apply the *Horner Scheme*. Let  $f(x) = \sum_{i=0}^n a_i x^i$  be such a polynomial approximation, then  $f(x)$  is efficiently evaluated as

$$f(x) = (\cdots((a_n) * x + a_{n-1}) * x \cdots + a_1) * x + a_0,$$

where to avoid a growth in operand lengths, roundings are needed in each cycle of the algorithm, i.e., after each multiply-add operation. But here the round-bits can easily be absorbed in a subsequent arithmetic operation, only at the very end a regular conversion may be needed, but normally the result is to be used in some calculation, hence again a conversion may be avoided.

For inner product calculations, the most accurate result is obtained if accumulation is performed in double precision, it will even be exact when performed in fixed-point arithmetic. However, if double precision is not available it is essential that a fast and optimal rounding is employed during accumulation of the product terms.

## VIII. CONCLUSIONS AND DISCUSSION

We have analyzed a general class of number representations for which truncation of a digit string yields the effect of rounding to nearest.

Concentrating on binary RN-represented operands, we have shown how a simple encoding, based on the ordinary 2's complement representation, allows trivial (constant time) conversion from 2's complement representation to the binary RN-representation. A simple parallel prefix (log time) algorithm is needed for conversion the other way. We have demonstrated how operands in this particular canonical encoding can be used at hardly any penalty in many standard calculations, e.g., addition and multiplication, with negation even being a constant time operation, which often simplifies the implementation of arithmetic algorithms.

The particular feature of the RN-representation, that rounding-to-nearest is obtained by truncation, implies that repeated roundings ending in some precision yields the same result, as if a single rounding to that precision was performed. In [Lee89] it was proposed to attach some

state information (2 bits) to a rounded result, allowing subsequent roundings to be performed in such a way, that multiple roundings yields the same result as a single rounding to the same precision. It was shown that this property holds for any specific IEEE-754 [IEE85] rounding mode, including in particular for the round-to-nearest-even mode. But these roundings may still require log-time incrementations, which is avoided with the proposed RN-representation.

The fixed point encoding immediately allows for the definition of corresponding floating point representations, which in a comparable hardware FPU implementations will be simpler and faster than those of equivalent IEEE standard conforming implementation.

Thus in applications where many roundings are needed, and conformance to the IEEE-754 standard is not required, employing the RN-representation it is possible to avoid the penalty of intermediate log-time roundings. Signal processing may be an application area where specialized hardware (ASIC or FPGA) is often used anyway, and the RN-representation can provide faster arithmetic with round to nearest operations at reduced area and delay.

## REFERENCES

- [BK82] R.P. Brent and H.T. Kung. A Regular Layout for Parallel Adders. *IEEE Transactions on Computers*, C-31(3):260–264, March 1982.
- [BM05] J.-L. Beuchat and J.-M. Muller. Multiplication Algorithms for Radix-2 RN-Codings and Two’s Complement Numbers. Technical report, INRIA, February 2005. Available at: <http://www.inria.fr/rrrt/rr-5511.html>.
- [Boo51] A.D. Booth. A Signed Binary Multiplication Technique. *Q. J. Mech. Appl. Math.*, 4:236–240, 1951.
- [IEE85] IEEE. *Std 754-1985 IEEE Standard for Binary Floating-Point Arithmetic*. Standards Committee of The IEEE Computer Society. 345 East 47th Street, New York, NY 10017, USA, 1985.
- [KM05] P. Kornerup and J.-M. Muller. RN-Coding of Numbers: Definition and some Properties. In *Proc. IMACS’2005*, July 2005. Paris.
- [Lee89] C. Lee. Multistep Gradual Rounding. *IEEE Transactions on Computers*, 38(4):595–600, April 1989.
- [NMLE00] A. Munk Nielsen, D.W. Matula, C.N. Lyu, and G. Even. An IEEE Compliant Floating-Point Adder that Conforms with the Pipelined Packet-Forwarding Paradigm. *IEEE Transactions on Computers*, 49(1):33–47, January 2000.