

On the design of IEEE compliant floating point units

Guy Even* and Wolfgang Paul
Univ. des Saarlandes
66123 Saarbrücken, Germany
E-mail:guy,wjp@cs.uni-sb.de

Abstract

Engineering design methodology recommends designing a system as follows: Start with an unambiguous specification, partition the system into blocks, specify the functionality of each block, design each block separately, and glue the blocks together. Verifying the correctness of an implementation reduces then to a local verification procedure.

We apply this methodology for designing a provably correct modular IEEE compliant Floating Point Unit. First, we provide a mathematical and hopefully unambiguous definition of the IEEE Standard which specifies the functionality. The design consists of: an adder, a multiplier, and a rounding unit, each of which is further partitioned. To the best of our knowledge, our design is the first publication that deals with detecting exceptions and trapped overflow and underflow exceptions as an integral part of the rounding unit in a Floating Point Unit. Our abstraction level avoids bit-level arguments while still enabling addressing crucial implementation issues such as delay and cost.

1. Introduction

Background. The IEEE Standard for Floating Point Arithmetic [1] defines the functionality of floating point arithmetic. Since its approval more than 10 years ago, the Standard has been an immense success, and all major Floating Point Units (FPU's) comply with it. Unfortunately, most designs are obtained by first designing an FPU that handles correctly most of the cases, and subsequent modifications are made to handle the remaining cases. The drawbacks of this approach are: (a) hard to predict performance and cost due to the modifications; (b) longer design process and complications in designs; and (c) hard to verify designs.

*Supported in part by Graduiertenkolleg "Effizienz und Komplexität von Algorithmen und Rechenanlagen", Universität des Saarlandes, and in part by the North Atlantic Treaty Organization under a grant awarded in 1996.

Previous work. There is a vast amount of literature on designing FPU's. The literature falls mainly into two categories: (a) Fundamental work that describes how FPU's are designed but are either not fully IEEE compliant or ignore details involved with the Standard such as exceptions. Among such works are [2, 3, 5, 6, 7]. (b) Implementation reports that describe existing FPU's while emphasizing features that enable improved performance.

Goals. We set forth the following goals concerning the design of an IEEE compliant FPU. (a) Provide an abstraction level that can bridge between a high level description and detailed issues involved in complying with the Standard. (b) Identify and characterize essential properties that explain the correctness of the circuit. (c) Formalize mathematical relations that can serve as "glue" for composing the building blocks together.

Overview and Contributions. This paper deals systematically with the design of IEEE compliant FPU's. A key issue is the choice of an abstraction level. We use an abstraction level suggested by Matula [4] which enables ignoring representation issues and simplifies greatly the description since one can avoid bit-level arguments. The Standard specifies the functionality but is hard to interpret. We believe that it is imperative to have a clear and structured specification. In particular, we define representable values, rounding, and exceptions. Modularity is obtained by defining a succinct property that the output of an adder and the input to the rounding unit should satisfy so that correctness is guaranteed. A rounding unit design is presented that detects exceptions correctly and deals correctly with trapped overflows and underflows. An addition algorithm is presented, and its correctness is discussed.

All the proofs are omitted and appear in the full version of the paper.

2. IEEE Standard - representable numbers

In this section we review which numbers are representable according to the IEEE Standard. The abstraction level ignores representation and uses factorings as suggested

by Matula [4]. A factoring factors a real number into three components: a sign factor, a scale factor, and a significand. We define the set of representable floating point numbers and describe their geometry.

Factorings

Every real number x can be factored into a sign factor (determined by a sign-bit), a scale factor (determined by an exponent), and a significand as follows:

$$x = (-1)^s \cdot 2^e \cdot f$$

The sign bit s is in $\{0, 1\}$, the exponent e is an integer, and the significand f is a non-negative real number. Usually the range of the significand is limited to the half open interval $[0, 2)$, but for example, intermediate results may have significands that are larger than 2. We henceforth refer to a triple (s, e, f) as a *factoring*, and $val(s, e, f)$ is defined to be $(-1)^s \cdot 2^e \cdot f$. A natural issue is that of unique representation. We postpone this issue until Sec. 3.1 in which we define a unique factoring called a *normalized factoring*.

We need to introduce a factoring of $\pm\infty$ as well. We do that by introducing a special exponent symbol e_∞ and a special significand symbol f_∞ that must appear together in a factoring (namely, an integer exponent cannot appear with f_∞ and a real significand cannot appear with e_∞). Thus, the factoring of ∞ is $(0, e_\infty, f_\infty)$, and the factoring of $-\infty$ is $(1, e_\infty, f_\infty)$.

Standard's representable real numbers

In this section we define which numbers are representable according to the IEEE Standard.

Every format (single, double, etc.) has two parameters attached to it: (a) n - the length of the exponent string; and (b) p - the length of the significand string (including the hidden bit).

The set of representable exponent values is the set of all integers between e_{\max} and e_{\min} . The values of e_{\max} and e_{\min} are determined by the parameter n as follows:

minimum exponent: $e_{\min} \triangleq 1 - bias$.

maximum exponent: $e_{\max} \triangleq 2^n - 2 - bias$.

where $bias \triangleq 2^{n-1} - 1$. This range of exponent values results from the biased binary representation of exponents.

The set of representable significand values is the set of all integral multiples of $2^{-(p-1)}$ in the half open interval $[0, 2)$. We distinguish between two ranges: representable significand values in the half open interval $[0, 1)$ are called *denormalized*, and representable significand values in the half open interval $[1, 2)$ are called *normalized*.

Not all combinations of representable exponent values and significand values are representable. Specifically, the denormalized significands can only appear with the exponent value e_{\min} . However, the normalized significands may appear with all the representable exponent values.

The Geometry of representable numbers

We depict the non-negative representable real numbers in Fig. 1; the picture for the negative representable numbers is symmetric.

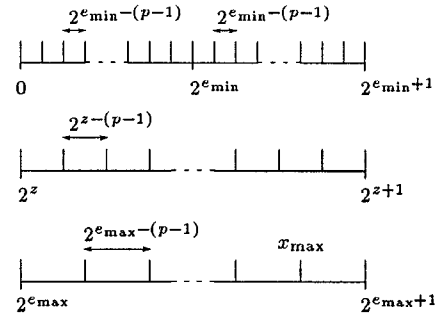


Figure 1. Geometry of representable numbers

The following properties characterize the representable numbers:

(1) For every exponent value z between e_{\min} and e_{\max} there are two intervals of representable numbers: $[2^z, 2^{z+1})$ and $(-2^{z+1}, -2^z]$. The gaps between consecutive representable numbers in these intervals are $2^{z-(p-1)}$.

(2) As the exponent value increases by one, the length of the interval $[2^z, 2^{z+1})$ doubles, and the gaps between the representable numbers double as well. Thus, the number of representable numbers per interval is fixed and it equals 2^{p-1} .

(3) The denormalized numbers, namely, the representable numbers in the interval $(-2^{e_{\min}}, 2^{e_{\min}})$, have an exponent value of e_{\min} and a significand value in the interval $[0, 1)$. The gaps between consecutive representable denormalized numbers are $2^{e_{\min}-(p-1)}$. Thus, the gaps in the interval $[0, 2^{e_{\min}})$ equal the gaps in the interval $[2^{e_{\min}}, 2^{e_{\min}+1})$. This property is called in the literature *gradual underflow* since the large gap between zero and $2^{e_{\min}}$ is filled with denormalized numbers.

3. Rounding - Definition

In this section we define rounding in round-to-nearest (even) mode.

3.1. Normalized factoring

Every real number has infinitely many factorings. We would like to obtain a unique factoring by requiring that $1 \leq f < 2$. However, there are two problems with this requirement. First, zero cannot be represented, and second, we do not want to have very small exponents. (We deal with large exponents later, during exponent-rounding). A normalized factoring is a unique factoring which avoids very small exponents and enables representing zero.

Definition 1 Normalization shift is the mapping η from the reals to factorings which maps every real x into $\eta(x) = (s, e, f)$ so that $x = \text{val}(\eta(x)) = \text{val}(s, e, f)$, i.e. value is preserved, and the following conditions hold:

1. If $\text{abs}(x) \geq 2^{e_{\min}}$, then $1 \leq f < 2$. In this case, we say that the significand f is a normalized significand.
2. If $\text{abs}(x) < 2^{e_{\min}}$, then $e = e_{\min}$ and $0 \leq f < 1$. In this case, we say that the significand f is a denormalized significand.

The sign-bit in the normalization shift of zero is not well defined. However, the Standard sets rules regarding the sign bit when the exact result equals zero. For example, $(+5) \cdot (-0) = -0$ and $5 - 5 = +0$. Therefore, we assume that the rounding unit is input the correct sign bit and the rounding unit does not change the sign bit. When $x = \pm\infty$, the normalization shift of x is $(s, e_{\infty}, f_{\infty})$, where $s = \text{sign}(x)$, and e_{∞}, f_{∞} denote the exponent and significand symbols used for factoring $\pm\infty$.

We also consider normalized factorings.

Definition 2 A factoring (s, e, f) is the normalized factoring of (s', e', f') (respectively x) if it satisfies $(s, e, f) = \eta(\text{val}(s', e', f'))$ (respectively $(s, e, f) = \eta(x)$). A factoring (s, e, f) is normalized if $(s, e, f) = \eta(\text{val}(s, e, f))$.

We use the term “normalized” for two different purposes: a normalized significand is a significand in the range $[1, 2)$, whereas a normalized factoring is a unique factoring representing a real number. The disadvantage of this terminology is that a normalized factoring can have a denormalized significand.

3.2. Significand rounding

Consider a significand $f \geq 0$. Suppose we would like to round f so that its binary representation has at most $p - 1$ bits to the right of the binary point (this implies p bits of precision when $f < 2$). We deal with rounding in round-to-nearest (even) mode.

First, we sandwich f by two consecutive integral multiples of $2^{-(p-1)}$ as follows:

$$q \cdot 2^{-(p-1)} \leq f < (q + 1) \cdot 2^{-(p-1)}, \quad (1)$$

Let q' be the even integer in $\{q, q + 1\}$. The significand rounding of f , denoted by $\text{sig_rnd}(f)$ is defined as follows: (a) If $q \cdot 2^{-(p-1)} \leq f < (q + 0.5) \cdot 2^{-(p-1)}$, then $\text{sig_rnd}(f) = q \cdot 2^{-(p-1)}$; (b) If $f = (q + 0.5) \cdot 2^{-(p-1)}$, then $\text{sig_rnd}(f) = q' \cdot 2^{-(p-1)}$; and (c) If $(q + 0.5) \cdot 2^{-(p-1)} < f < (q + 1) \cdot 2^{-(p-1)}$, then $\text{sig_rnd}(f) = q \cdot 2^{-(p-1)}$.

Define the signal SIG_INEXACT as follows:

$$\text{SIG_INEXACT} = 1 \quad \text{if } \text{sig_rnd}(f) \neq f$$

The SIG_INEXACT signal is required for detecting the loss of accuracy.

Note that the rounding modes: round to $+\infty$ and round to $-\infty$, as defined in the Standard, are non-symmetric rounding modes. Namely, in these rounding modes, significand rounding depends also on the sign bit. For simplicity, we omit the sign bit as an argument of significand rounding. However, in our rounding unit depicted in Fig. 2, the sign bit is input to the significand rounding box.

3.3. Post-Normalization

It is important to note that if $f \in [0, 2)$, then the result of significand-rounding, $\text{sig_rnd}(f)$, is in the range $[0, 2]$. The case that $\text{sig_rnd}(f) = 2$ is called *significand overflow*. When significand overflow occurs, the significand is set to 1 and the exponent is incremented. This normalization shift takes place only when a significand overflow occurs. Therefore, post-normalization, denoted by $\text{post_norm}(s, e, f)$, is defined as follows:

$$\text{post_norm}(s, e, f) \triangleq \begin{cases} (s, e + 1, 1) & \text{if } f = 2 \\ (s, e, f) & \text{otherwise} \end{cases}$$

3.4. Exponent rounding

Exponent rounding maps factorings into factorings and deals with the case that the absolute value of a factoring is too large. According to the Standard, exponent rounding is bypassed when trapped overflow exception occurs. We define only exponent rounding in round-to-nearest (even) mode.

Recall that e_{∞} and f_{∞} denote the exponent and significand symbols used for factoring $\pm\infty$. Let $x_{\max}^* = 2^{e_{\max}} \cdot (2 - 2^{-p})$. The exponent-rounding of (s, e, f) is defined as follows:

$$\text{exp_rnd}(s, e, f) \triangleq \begin{cases} (s, e_{\infty}, f_{\infty}) & \text{if } 2^e \cdot f \geq x_{\max}^* \\ \eta(\text{val}(s, e, f)) & \text{otherwise} \end{cases}$$

The motivation for the definition of the threshold of x_{\max}^* is as follows. Let x_{\max} denote the largest representable number, namely, $x_{\max} = 2^{e_{\max}} \cdot (2 - 2^{-(p-1)})$. If we were not limited by e_{\max} , then the next representable number

would be $2^{e_{\max}} \cdot 2$. The midpoint between these numbers is exactly x_{\max}^* . The binary representation of the significand of x_{\max} is $1.11 \dots 1$. Therefore, the significand is an odd integral multiple of $2^{-(p-1)}$. Hence, x_{\max}^* would be rounded up in round-to-nearest (even) mode. This is why numbers greater than or equal to x_{\max}^* are rounded to infinity. A similar argument holds for rounding of values not greater than $-x_{\max}^*$ to minus infinity.

Note that if $1 \leq f < 2$ is an integral multiple of $2^{-(p-1)}$, then we can simplify the definition of exponent rounding, and round to infinity if $e > e_{\max}$.

3.5. Rounding

Rounding is a mapping of reals into factorings. The rounding of x , denoted by $r(x)$, is defined as follows: Let (s, e, f) denote the normalized factoring of x , i.e. $(s, e, f) = \eta(x)$, then

$$r(x) \triangleq \text{exp_rnd}(\text{post_norm}(s, e, \text{sig_rnd}(f))) \quad (2)$$

Therefore, rounding is the composition of four functions: (a) normalization shift - to obtain the right factoring; (b) significand rounding - to obtain the right precision (namely, limit the length of the significand); (c) post-normalization - to correct the factoring in case significand overflow occurs; and (d) exponent rounding - to limit the range of the value of the factoring. Note that $r(x)$ is a normalized factoring.

4. Standard's exceptions

Five exceptions are defined by the Standard: invalid operation, division by zero, overflow, underflow, and inexact. The first two exceptions have to do only with the strings given as input to an operation rather than with the outcome of a computation. Hence, we focus on the last three exceptions. The main advantage is that we can use our notation for short and precise definitions of the exceptions.

Exception handling consists of two parts: signaling the occurrence of an exception by setting a status flag and invoking a trap handler (if it is enabled). We start by defining each exception, namely, we define the conditions that cause an exception to occur.

Assumptions: If one of the operands is infinite or not-a-number, then the overflow, underflow and inexact exceptions do not occur. Hence, we assume both operands represent finite values. Moreover, we assume that the exact result is also finite. (An infinite exact result generates a division by zero exception, an invalid exception, or does not signal any exception).

4.1. Definitions

Defining exceptions requires two additional definitions.

Definition 3 A normalization shift with unbounded exponent range is the mapping $\hat{\eta}$ from the non-zero reals to factorings which maps every real x into a factoring (s, \hat{e}, \hat{f}) so that $x = \text{val}(\hat{\eta}(x))$, i.e. value is preserved, and $\hat{f} \in [1, 2)$.

In other words, $\hat{\eta}(x)$ is the normalized factoring of $x \neq 0$ if there were no lower bound on the exponent range.

Definition 4 The rounding with an unbounded exponent range is the mapping \hat{r} of non-zero real into factorings defined by:

$$\hat{r}(x) \triangleq \text{post_norm}((s, \hat{e}, \text{sig_rnd}(\hat{f}))), \text{ where } \hat{\eta}(x) = (s, \hat{e}, \hat{f})$$

Notation: Throughout this section we use x to denote the exact (finite) value of the result of an operation. Let $\eta(x) = (s, e, f)$, and $\hat{\eta}(x) = (s, \hat{e}, \hat{f})$.

Overflow. Informally, overflow occurs when the magnitude of the result is larger than x_{\max} . However, this is not precise, because significand rounding can lower the magnitude back into the range of representable numbers. The precise definition is given below.

Definition 5 Let x denote the exact result, and x_{\max} the largest representable number. An overflow exception occurs if

$$\text{val}(\hat{r}(x)) > x_{\max} \text{ or } \text{val}(\hat{r}(x)) < -x_{\max}$$

Note that the definition of overflow is with respect to rounding with an unbounded exponent range.

The following claim facilitates the detection of overflow, since we do not compute $\hat{r}(x)$.

Claim 1 Let (s, e, f) denote the normalized factoring of the exact result. An overflow exception occurs iff $2^e \cdot \text{sig_rnd}(f) > x_{\max}$ or, equivalently, iff $e > e_{\max}$ or $(e = e_{\max} \text{ and } \text{sig_rnd}(f) = 2)$.

Underflow. The definition of underflow in the Standard is extremely complicated. In the definitions, we follow the language of the Standard, although some of the definitions include irrelevant cases. Informally, underflow occurs when two conditions occur: (a) *tininess* - The magnitude of the result is below $2^{e_{\min}}$; and (b) *loss-of-accuracy* - Accuracy is lost when the tiny result is represented with a denormalized significand. The Standard gives two definitions for each of these conditions, and thus, the Standard gives four definitions of underflow, each of which conforms with the Standard. Note that the same definition must be used by an implementation for all operations. However, the Standard does not state that the same definition should be used for all precisions.

The two definitions of tininess, *tiny-after-rounding* and *tiny-before-rounding*, are defined below.

Definition 6 *tiny-after-rounding occurs if $x \neq 0$ and*

$$0 < \text{abs}(\widehat{r}(x)) = 2^{\widehat{e}} \cdot \text{sig_rnd}(\widehat{f}) < 2^{e_{\min}}$$

tiny-before-rounding occurs if

$$0 < \text{abs}(x) = 2^e \cdot f < 2^{e_{\min}}$$

In round-to-nearest (even) mode, *tiny-after-rounding* occurs iff $2^{e_{\min}} \cdot 2^{-p-1} < \text{abs}(x) < 2^{e_{\min}} \cdot (1 - 2^{-p-1})$. In general, *tiny-after-rounding* implies *tiny-before-rounding*.¹

The two definitions of loss of accuracy are defined below.

Definition 7 *loss-of-accuracy-a (also called denormalization loss) occurs if $x \neq 0$ and*

$$r(x) \neq \widehat{r}(x)$$

loss-of-accuracy-b (also called inexact result) occurs if

$$\text{val}(r(x)) \neq x$$

The Standard is unclear about whether loss-of-accuracy should also occur when overflow occurs. We have chosen to define loss-of-accuracy so that the two types of loss-of-accuracy occur when overflow occurs. However, we are interested in detecting loss-of-accuracy when tininess occurs, and not when overflow occurs. In this context, *loss-of-accuracy-b* simply means that significand rounding introduces an error. The meaning of *loss-of-accuracy-a* is more complicated: When f is normalized (namely, in the range $[1, 2)$), then $\eta(x) = \widehat{\eta}(x)$, and if no overflow occurs, then *loss-of-accuracy-a* does not occur. However, when f is denormalized, then $e = e_{\min}$ and the binary representation of f has $e - \widehat{e}$ leading zeros. Therefore, $\text{sig_rnd}(\widehat{f})$ is in general different from $2^{e - \widehat{e}} \cdot \text{sig_rnd}(f)$, because extra nonzero digits are “shifted in”. Note that *loss-of-accuracy-a* implies *loss-of-accuracy-b*.

The definition of underflow depends on whether the underflow trap handler is enabled or disabled.

Definition 8 *If the underflow trap handler is disabled, then underflow occurs if both tininess and loss-of-accuracy occur. If the underflow trap handler is enabled, then underflow occurs if tininess occurs.*

Note that tininess in Def. 8 means either *tiny-after-rounding* or *tiny-before-rounding*. Similarly, loss-of-accuracy means either *loss-of-accuracy-a* or *loss-of-accuracy-b*.

Inexact. The Standard defines the inexact exception twice, and the two definitions turn out to be equivalent. The first definition says that an inexact exception occurs when the value of the delivered result does not equal the exact result. This has already been defined in Def. 7 as *loss-of-accuracy-b*. The second definition is as follows.

¹Note that $0 < 2^e \cdot \text{sig_rnd}(f) < 2^{e_{\min}}$ in round-to-nearest (even) mode is equivalent to $2^{e_{\min}} \cdot 2^{-p} < \text{abs}(x) < 2^{e_{\min}} \cdot (1 - 2^{-p})$.

Definition 9 *An inexact exception occurs if $\text{val}(r(x)) \neq x$ or if an overflow occurs with the trap handler disabled.*

Since overflow implies *loss-of-accuracy-b*, it follows that both definitions are equivalent.

The following claim facilitates the detection of an inexact exception.

Claim 2 $\text{val}(r(x)) = x$ iff $\text{sig_rnd}(f) = f$ and $\text{exp_rnd}(s, e, \text{sig_rnd}(f)) = (s, e, \text{sig_rnd}(f))$.

Note that the overflow and underflow trap handlers have precedence over the inexact trap handler.

4.2. Trapped overflows and underflows

The Standard sets special rules for the delivered result when an overflow or an underflow exception occurs and the corresponding trap handler is enabled. We call these cases *trapped overflows* and *trapped underflows*, respectively. In these cases, the exponent is “wrapped” so that the delivered result is brought back into the range of representable normalized numbers. We formalize the Standard’s requirements below.

Definition 10 *Exponent wrapping in trapped overflow and underflow exceptions:*

1. *If a trapped overflow occurs, then the value of the delivered result is $\text{val}(r(x \cdot 2^{-\alpha}))$, where $\alpha = 3 \cdot 2^{n-2}$.*
2. *If a trapped underflow occurs, then the value of the delivered result is $\text{val}(r(x \cdot 2^{\alpha}))$.*

The difficulty that this “wrapped exponent” requirement causes is that we do not have the exact result x , but only the factoring $(s, e, \text{sig_rnd}(f))$. Our goal is to integrate these requirements with the computation of rounding. The following claim shows that this can be easily done in the case of overflow.

Claim 3 *Let $(s, e', f') = \text{post_norm}(s, e, \text{sig_rnd}(f))$, where $\eta(x) = (s, e, f)$. If a trapped overflow occurs, then*

$$r(2^{-\alpha} \cdot x) = (s, e' - \alpha, f')$$

The implication of Claim 3 is that when a trapped overflow occurs, the only fix required is to subtract α from the exponent.

The analogous claim for the case of underflow does not hold. The reason being that $(s, e + \alpha, f)$ is not the normalized factoring of $2^{\alpha} \cdot x$. Recall, that f is denormalized, and hence, its binary representation has leading zeros. However, the normalized factoring of $2^{\alpha} \cdot x$ has a normalized significand. Therefore, multiplying x by 2^{α} effects both the exponent and the significand in the normalized factoring.

Dealing with a trapped underflow exception is simple in addition and subtraction operations since the rounded result is exact (see Sec. 7.3). Therefore, $\text{val}(r(2^\alpha \cdot x)) = \text{val}(s, e + \alpha, f)$, and hence, one only needs to normalize f and update the exponent by adding α minus the number of leading zeros in f .

Dealing with a trapped underflow exception in multiplication relies on the availability of $\hat{\eta}(x)$. Consider the factoring $(s, \hat{e}, \hat{f}) = \hat{\eta}(x)$. Note that since underflow occurs, it follows that $\hat{e} < e_{\min}$. We formalize an analogous claim for trapped underflow exceptions using $\hat{\eta}(x)$.

Claim 4 *Let $(s, \hat{e}, \hat{f}) = \hat{\eta}(x)$, and $(s, e', f') = \text{post_norm}(s, \hat{e}, \text{sig_rnd}(\hat{f}))$. If a trapped underflow occurs, then*

$$r(2^\alpha \cdot x) = (s, e' + \alpha, f')$$

5. Rounding - Computation

In this section we deal with the issue of how to design a rounding unit. Our goal is to design a rounding unit that also detects the exceptions: overflow, underflow, and inexact. Moreover, the rounding unit wraps the exponent in case of a trapped overflow or underflow exception. Significant rounding is presented using representatives that require less bits of precision while guaranteeing correct rounding. We relate representatives with sticky bits, and present a block diagram of a rounding unit capable of detecting and dealing with exceptions.

5.1. Representatives

Definition 11 *Suppose α is an integer. Two real numbers x_1 and x_2 are α -equivalent, denoted by $x_1 \stackrel{\alpha}{\cong} x_2$ if there exists an integer q such that $x_1, x_2 \in (q2^{-\alpha}, q2^{-\alpha} + 2^{-\alpha})$ or $x_1 = x_2 = q2^{-\alpha}$.*

The binary representations of α -equivalent reals must agree in the first α positions to the right of the binary point. Note that in case a number has two binary representations, we choose the finite representation (for example, 0.1 rather than 0.0111...).

We choose α -representatives of the equivalence classes as follows.

Definition 12 *Let x denote a real number and α an integer. Let q denote the integer satisfying: $q2^{-\alpha} \leq x < (q+1)2^{-\alpha}$. The α -representative of x , denoted by $\text{rep}_\alpha(x)$, is defined as follows:*

$$\text{rep}_\alpha(x) \triangleq \begin{cases} q2^{-\alpha} & \text{if } x = q2^{-\alpha} \\ (q + 0.5) \cdot 2^{-\alpha} & \text{if } x \in (q2^{-\alpha}, (q + 1)2^{-\alpha}) \end{cases}$$

Note that for a fixed α , the set of all α -representatives equals the set of all integral multiples of $2^{-\alpha-1}$. Thus, the α -representatives have $\alpha + 1$ bits of precision beyond the binary point. Moreover, the least significant bit of the binary representation of an α -representative can be regarded as a flag indicating whether the corresponding equivalence class is a single point or an open interval.

The following claim summarizes the properties of representatives that we use in significant rounding.

Claim 5 *Let $f' = \text{rep}_p(f)$, then*

$$\begin{aligned} \text{sig_rnd}(f) &= \text{sig_rnd}(f') \\ \text{sig_rnd}(f) = f &\text{ iff } \text{sig_rnd}(f') = f' \\ \text{sig_rnd}(f') = f' &\text{ iff } f' = q \cdot 2^{-(p-1)} \end{aligned}$$

Claim 5 shows that one can substitute f with $\text{rep}_p(f)$ for the purpose of computing $\text{sig_rnd}(f)$. Not only do they generate the same significant rounding, but they also generate the same SIG_INEXACT signal. Moreover, the SIG_INEXACT signal can be easily computed from $\text{rep}_p(f)$ because SIG_INEXACT = 1 iff $\text{rep}_p(f)$ is not an integral multiple of $2^{-(p-1)}$. Since $\text{rep}_p(f)$ is an integral multiple of $2^{-(p+1)}$, it follows that SIG_INEXACT is the OR of the two least significant bits of the binary representation of $\text{rep}_p(f)$.

The following claim shows a simple and sufficient condition for two reals to have the same rounded factoring. We heavily rely on this claim for gluing the rounding unit with the adder.

Claim 6 *If $\eta(x) = (s_x, e_x, f_x)$, $\eta(y) = (s_y, e_y, f_y)$ and $x \stackrel{p-e_x}{\cong} y$, then (a) $s_x = s_y$, $e_x = e_y$ and $f_x \stackrel{p}{\cong} f_y$; and (b) $r(x) = r(y)$.*

Consider an exact result x and a computed result y . Claims 5 and 6 imply that the inexact exception is detected correctly even if we use y rather than x .

Based on Claim 2 we can detect an inexact exception as follows:

Corollary 7 *An inexact exception occurs iff an untrapped overflow exception occurs or $\text{rep}_p(f)$ is not an integral multiple of $2^{-(p-1)}$.*

5.2. Sticky-bits and representatives

The common way of “shortening” a significant f whose binary representation requires more than $p - 1$ bits to the right of the binary point is called a *sticky-bit computation*. In a sticky-bit computation, one replaces the tail of the binary representation of f starting in position $p + 1$ to the right of the binary point with the OR of the bits in this tail. We show that the value of the binary string obtained by this substitution equals the p -representative of f . Thus our use

of representatives captures the common practice of using a sticky-bit.

Claim 8 Let f denote a significand and let $\text{bin}(f) = f_0.f_1f_2\dots$ denote the binary representation of f , namely, $f = \sum_i f_i \cdot 2^{-i}$. Define:

$$\begin{aligned} \text{sticky_bit}(f, p) &\triangleq \text{OR}(f_{p+1}, f_{p+2}, \dots) \\ \text{sticky}(f, p) &\triangleq \sum_{i \leq p} f_i \cdot 2^{-i} + \text{sticky_bit}(f, p) \cdot 2^{-p-1} \end{aligned}$$

Then,

$$\text{rep}_p(f) = \text{sticky}(f, p)$$

5.3. The significand-rounding decision

Given the representative of $f \in [0, 2)$, significand rounding reduces to a decision whether the p most significant bits of the binary representation of $\text{rep}_p(f)$ should be incremented or not. The decision depends on: the three least significant bits of the binary representation of $\text{rep}_p(f)$, the sign bit, and the rounding mode.

The three least significant bits of the binary representation of $\text{rep}_p(f)$ are called: the LSB (which is $p-1$ positions to the right of the binary point), the round-bit, and the sticky-bit. Using the notation used in defining significand rounding: the LSB is needed to determine the parity of q , whereas the round-bit and sticky-bit determine which of the three cases is selected.

5.4. Design - Block Diagram

In this section we describe a block-diagram of a rounding unit. Obviously other more efficient implementations are possible and have been implemented. The advantages of our design are: (a) Well defined requirements from the inputs to facilitate the usage of the rounding unit for different operations. (b) Correct and full detection of exceptions. (c) Correct handling of trapped underflow and overflow exceptions. Whenever no confusion is caused, we refer to a representation of a number (a significand, exponent, etc.) as the number itself. Fig. 2 depicts a block diagram of a rounding unit.

Input. The inputs consist of: a factoring (s_{in}, e_{in}, f_{in}) , the rounding mode, flags called UNF_EN and OVF_EN indicating whether the underflow trap and overflow trap handlers respectively are enabled.

Assumptions. We assume that the inputs satisfy the following two assumptions:

1. Let x denote the exact result. For example, x is the exact sum of two values that are added by the adder

and rounded by the rounding unit. Let $(s, \widehat{e}_x, \widehat{f}_x) = \widehat{\eta}(x)$. Then,

$$\text{val}(s_{in}, e_{in}, f_{in}) \stackrel{p-\widehat{e}_x}{=} x$$

2. $\text{abs}(x)$ is not too large or small so that wrapping the exponent when a trapped underflow or overflow occurs produces a number within the range of numbers representable with normalized significands. Namely,

$$2^{e_{\min}-\alpha} \leq \text{abs}(\widehat{r}(x)) < 2^{e_{\max}+1+\alpha}$$

Output. The rounding unit outputs a factoring $(s_{out}, e_{out}, f_{out})$ that equals: (a) $r(x)$ if no trapped overflow or underflow occurs; (b) $r(x \cdot 2^{-\alpha})$ if a trapped overflow occurs; and (c) $r(x \cdot 2^{\alpha})$ if a trapped underflow occurs. According to Claims 3, 4, and 6, this specification can be reformulated as (a) $r(\text{val}(s_{in}, e_{in}, f_{in}))$ if no trapped overflow or underflow occurs; (b) $r(\text{val}(s_{in}, e_{in} - \alpha, f_{in}))$ if a trapped overflow occurs; and (c) $r(\text{val}(s_{in}, e_{in} + \alpha, f_{in}))$ if a trapped underflow occurs.

Additional flags are output: OVERFLOW, TINY and SIG_INEXACT. The OVERFLOW flag signals whether an overflow exception occurred. The TINY flag signals whether *tiny-before-rounding* occurred. The SIG_INEXACT signals whether significand rounding introduces a rounding error, and is used for detecting an inexact exception occurred (according to Coro. 7, inexact occurs iff SIG_INEXACT OR (OVERFLOW AND $\overline{\text{OVF_EN}}$)).

Functionality. We describe the functionality of each block in Fig.2:

1. The normalization shift box deals primarily with computing $\eta(\text{val}(s_{in}, e_{in}, f_{in}))$. However, it also deals with trapped underflows and partially deals with trapped overflows. Defining the functionality of the normalization shift box requires some notation: Let $(s_{in}, e'_{in}, f'_{in}) = \eta(\text{val}(s_{in}, e_{in}, f_{in}))$ and let $(s_{in}, \widehat{e}_{in}, \widehat{f}_{in}) = \widehat{\eta}(\text{val}(s_{in}, f_{in}, e_{in}))$. Then TINY, OVF_1 , e^n , f^n are defined as follows:

$$\text{TINY} = 1 \quad \text{if } 0 < 2^{e_{in}} \cdot f_{in} < 2^{e_{\min}}$$

$$\text{OVF}_1 = 1 \quad \text{if } 2^{e_{in}} \cdot f_{in} \geq 2^{e_{\max}+1}$$

$$(e^n, f^n) = \begin{cases} (e'_{in} - \alpha, f'_{in}) & \text{if } \text{OVF_EN} \text{ and } \text{OVF}_1 \\ (e_{in} + \alpha, f_{in}) & \text{if } \text{UNF_EN} \text{ and } \text{TINY} \\ (e'_{in}, f'_{in}) & \text{otherwise} \end{cases}$$

Note that when OVF_1 and $\overline{\text{OVF_EN}}$, the values of e^n and f^n are not important because e_{out} and f_{out} are fully determined in the exponent rounding box according to the rounding mode and the sign bit.

2. The $rep_p()$ box computes $f^1 = rep_p(f^n)$, meaning that the final sticky-bit is computed. Typically, a “massive” sticky-bit computation takes place before the rounding unit, and here, only the round-bit and sticky-bit are ORed, if the guard-bit turns out to be the LSB.
3. The significand rounding outputs $f^2 = sig_rnd(f^1)$ and two flags: SIG-OVF and SIG_INEXACT. The SIG-OVF flag signals the case that $f^2 = 2$, known as significand overflow. The SIG_INEXACT flag signals whether $f^2 \neq f^1$.
4. The post-normalization box uses the SIG-OVF flag as follows: If $f^2 = 2$, then the exponent needs to be incremented, and the significand needs to be replaced by 1. This can be obtained by ORing the two most-significant bits of the binary representation of f^2 , which is depicted in the figure by an OR gate that is input these two most-significant bits. Note, that since $f^2 \in [0, 2]$, the binary representation of f^2 has two bits to the left of the binary point.
5. The exponent adjust box deals with two issues caused by significand overflow: (a) a significand overflow might cause the rounded result to overflow; and (b) a significand overflow might cause a denormalized significand f^n to be replaced with a normalized significand $f^3 = 1$. (We later show that case (b) does not occur in addition/subtraction.)
 The output OVF_2 is set to one if $e^2 = e_{max} + 1$. The output e^3 is defined as follows: (a) If OVF_EN and OVF_2 , namely, significand rounding caused a trapped overflow, then $e^3 = e_{max} + 1 - \alpha$. (b) If $msb(f^3) = 1$ and TINY, namely, significand rounding causes a tiny result to have a normalized significand, then e^3 is set to the “normalized” representation of e_{min} . Note that this case deals with the representation of the exponent, namely, whether the representation reserved for denormalized significands should be used or the representation reserved for normalized significands. (c) Otherwise, $e^3 = e^2$. Note that if an untrapped overflow exception occurs, then the value of e^3 is not important because it is overridden by the exponent rounding.
6. The exponent rounding box performs exponent rounding outputs an exponent and significand (e_{out}, f_{out}) that equal: (a) (e^3, f^3) if $\overline{OVERFLOW}$ and $\overline{OVF_EN}$; (b) (e_{max}, f_{max}) if $\overline{OVERFLOW}$ and $\overline{OVF_EN}$ and round to x_{max} ; and (c) (e_{∞}, f_{∞}) if $\overline{OVERFLOW}$ and $\overline{OVF_EN}$ and round to ∞ . The decision of whether an overflowed result is rounded to x_{max} or to ∞ is determined by the rounding mode and the sign bit.

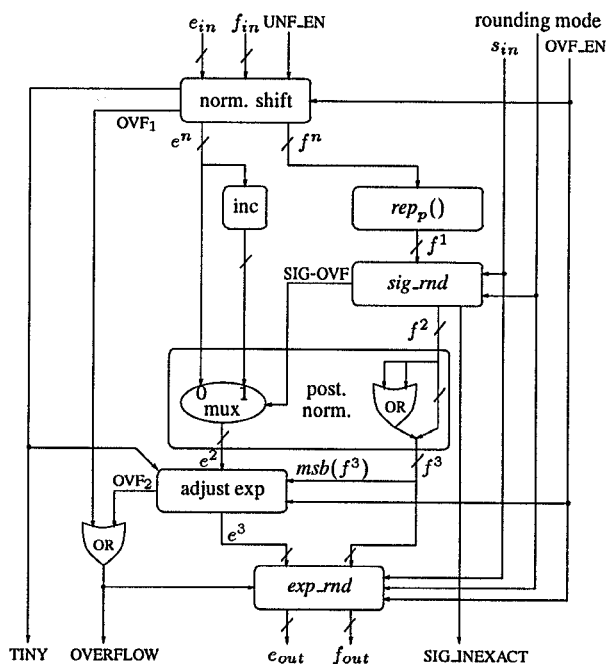


Figure 2. Rounding unit

6. Gluing rounding with other operations

In this section we discuss how the rounding unit is glued with an adder or a multiplier (or other functional units). We provide some mathematical tools that facilitate this task.

6.1. General framework

How does one design a circuit that meets the Standard’s requirement that “each of the operations shall be performed as if it first produced an intermediate result correct to infinite precision and with unbounded range, and then coerced this intermediate result to fit the destination’s format”? Obviously, one does not perform exact computations because that would be too expensive and slow. Fig. 3 depicts how infinite precision calculation is avoided. The functional unit is typically a fixed point unit, for example, an adder or a multiplier. Cost and delay is saved by minimizing the precision of the functional unit. Given two operands, rather than computing $op(a, b)$ with infinite precision, we pre-process the operands to bound the precision, and operate on a' and b' . The bounded precision result, denoted by $op(a', b')$, satisfies $r(op(a, b)) = r(op(a', b'))$, and thus, correctness is obtained. We guarantee that $r(op(a, b)) = r(op(a', b'))$ by using Claim 6, namely, by insuring that $op(a', b')$ and $op(a, b)$ are $(p - e)$ -equivalent. (When a trapped underflow occurs,

we need more precision, namely, $op(a', b') \stackrel{p-\hat{e}}{\equiv} op(a, b)$, where \hat{e} is the exponent in $\hat{\eta}(x)$. This succinct condition greatly simplifies the task of correctly gluing together the functional unit and the rounding unit.

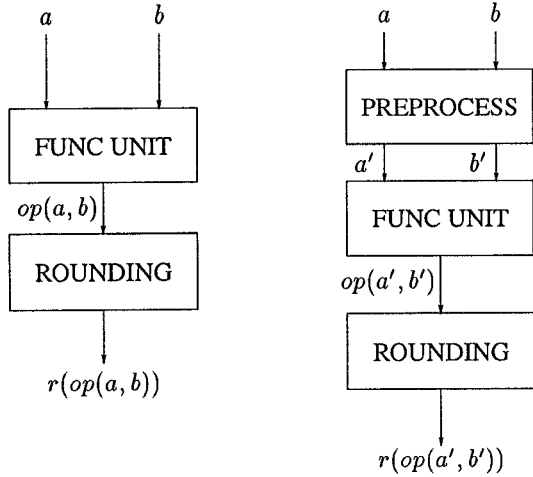


Figure 3. Avoiding infinite precision

6.2. Mathematical tools

The following claim is used in proving the correctness of rounding, addition, detection of exceptions, and handling of trapped overflows and underflows. It enables a symbolic treatment of the manipulations performed during various operations without having to deal with the peculiarities of representation.

Claim 9 *Let x, y denote real numbers such that $x \stackrel{\alpha}{\equiv} y$, for an integer α . Let $\alpha' \leq \alpha$ denote an integer. Then the following properties hold:*

1. $rep_{\alpha}(-x) = -rep_{\alpha}(x)$, and hence, $-x \stackrel{\alpha}{\equiv} -y$.
2. $x \stackrel{\alpha'}{\equiv} y$.
3. For every integer i : $x2^{-i} \stackrel{\alpha+i}{\equiv} y2^{-i}$.
4. $x + q2^{-\alpha'} \stackrel{\alpha}{\equiv} y + q2^{-\alpha'}$, for every integer q .

Claim 9 should be interpreted as follows: (1) $rep_{\alpha}(x)$ can be computed by computing $rep_{\alpha}(abs(x))$ and multiplying it by $sign(x)$. (2) Recall that an α -representative has $\alpha + 1$ bits of precision beyond the binary point. Thus, having too many bits of precision does not disturb correct rounding. (3) If two reals have identical binary representation up to bit position α beyond the binary point, then after division by 2^i (shift right by i positions) they have identical binary

representations up to bit position $\alpha + i$. Conversely, a shift to the left reduces the precision beyond the binary point. (4) Adding an integral multiple of $2^{-\alpha}$ to two reals whose binary representations agree up to bit position α beyond the binary point, does not create a disagreement in these bit positions.

7. Addition

In this section we present an addition algorithm and prove its correctness. By correctness we refer to delivering the correct result and detecting the exceptions. In particular, we provide a proof that 3 additional bits of precision suffice for correct rounding. We deal here only with adding finite valued factorings and ignore the Standard's requirements concerning the sign-bit when the exact result equals zero. The special cases, such as NAN's, infinite valued factorings, and the sign of zero results are handled by additional logic which we do not specify. Whenever no confusion is caused, we refer to a representation of a number (a significand, exponent, etc.) as the number itself.

7.1. Addition algorithm

The input to the adder consists of two factorings (s_1, e_1, f_1) and (s_2, e_2, f_2) which are normalized factorings. The adder outputs the factoring (s', e_1, g') , and this factoring is rounded by the rounding unit which is described in Sec. 5.4. The algorithm is divided into three stages: Preprocessing, Adding, and Rounding as follows:

Preprocessing

1. **Swap.** Swap the operands, if necessary, so that $e_1 \geq e_2$. This step is computed by comparing the exponents and selecting the operands accordingly. For simplicity, we assume that $e_1 \geq e_2$.
2. **Alignment Shift.** Replaces f_2 with f'_2 , where $f'_2 \triangleq f_2 \cdot 2^{-\delta}$ and $\delta = \min\{e_1 - e_2, p + 2\}$. This step is computed by shifting the binary representation of f_2 by δ positions to the right.
3. **Compute Representative.** Compute $f''_2 \triangleq rep_{p+1}(f'_2)$. This computation is, in effect, a sticky-bit computation as described in Claim 8. Note that the binary representation of f''_2 has $p + 2$ bits to the right of the binary point (namely, exactly 3 bits of precision are added to the significands).

Functional Unit

Add Significands. Output the factoring is (s', e_1, g') , where s' and g' are the sign and magnitude of $(-1)^{s_1} \cdot f_1 + (-1)^{s_2}$.

f_2'' . Note that g' has $p + 2$ bits to the right of the binary point and 2 bits to the left of the binary point.

Rounding

The Rounding unit receives the factoring (s', e_1, g') and additional inputs as described in Sec. 5.4.

Note, that the alignment shift and the representation computation can be performed in parallel since the shift amount $\delta = e_1 - e_2$ can be computed during the swap stage. Moreover, the width of the shifter equals $p + 2$, namely, the bits that are shifted beyond position $p + 1$ to the right of the binary point are discarded by the shifter. These discarded bits participate in the sticky-bit computation.

7.2. Correctness

In this section we discuss the correctness of the addition algorithm.

Let $x = \text{val}(s_1, e_1, f_1)$, $y = \text{val}(s_2, e_2, f_2)$, and $\hat{\eta}(x + y) = (s, \hat{e}, \hat{f})$. Let $y'' = \text{val}(s_2, e_1, f_2'')$. Since the factoring (s', e_1, g') is obtained by adding factorings that represent x and y'' , it follows that $\text{val}(s', e_1, g') = x + y''$. Based on the correctness of the rounding-unit, proving correctness reduces to proving the following claim.

Claim 10 $x + y \stackrel{p-\hat{e}}{=} x + y''$

Note that we need all the 3 additional bits of precision only if $e = e_1 - 1$. In this case, the normalization shift shifts the significand by one position to the left, and only 2 additional bits of precision are left, without them correct rounding cannot be performed.

7.3. Remarks

In this section we point out properties specific to addition and subtraction.

The following claim deals with the underflow exception in the addition operation. In particular, it shows that all four definitions of the underflow exception are identical with respect to addition and that a trapped underflow never occurs.

Claim 11 *Suppose two finite valued factorings are added, then:*

1. *tiny-after-rounding occurs iff tiny-before-rounding occurs.*
2. *If tiny-before-rounding (or tiny-after-rounding) occurs, then the adder outputs the exact sum, namely, neither loss-of-accuracy-a nor loss-of-accuracy-b occurs.*

Therefore, during addition, if the underflow trap handler is disabled, then the underflow exception never occurs. If the underflow trap handler is enabled, the underflow exception occurs iff $0 < f^n < 1$ where f^n denotes the significand output by the normalization shift box in the rounding unit.

An interesting consequence of Claim 11 is that the exponent adjust box in the rounding unit can be simplified for rounding sums and differences. Since a denormalized result is always precise, significand rounding cannot cause a denormalized significand to become normalized (thanks to Marc Daumas for pointing this out).

8. Multiplication

Multiplication is straightforward. The input consists of two factorings (s_1, e_1, f_1) and (s_2, e_2, f_2) which are normalized factorings. The multiplier outputs the factoring (s', e', f') , and this factoring is rounded by the rounding unit which is described in Sec. 5.4. The outputs are defined by $s' = \text{XOR}(s_1, s_2)$, $e' = e_1 + e_2$, and $f' = f_1 \cdot f_2$. Note, that if $f_1, f_2 \geq 1$, then one can avoid having a double length significand f' by computing $f' = \text{rep}_p(f_1 \cdot f_2)$.

Acknowledgments

We would like to thank Marc Daumas, Arno Formella, John Hauser, David Matula, Silvia Müller, Asger Nielsen, and Thomas Walle for many helpful suggestions and discussions.

References

- [1] "IEEE Standard for Binary Floating-Point Arithmetic", ANSI/IEEE Standard 754-1985.
- [2] J. T. Coonen, "Specification for a Proposed Standard for Floating Point Arithmetic", Memorandum ERL M78/72, University of California, Berkeley, 1978.
- [3] Isreal Koren, *Computer Arithmetic Algorithms*, Prentice-Hall, 1993.
- [4] David Matula, "Floating Point Representation", manuscript, May 1996.
- [5] Amos R. Omondi, *Computer Arithmetic Systems: Algorithms, Architecture and Implementations*, Prentice Hall, 1994.
- [6] Uwe Sparmann, *Strukturbasierte Testmethoden für arithmetische Schaltkreise*, Ph.D. Dissertation, Universität des Saarlandes, 1991.
- [7] Shlomo Waser and Michael J. Flynn, *Introduction to Arithmetic for Digital Systems Designers*, Holt, Reiner & Winston, 1982.