

# A MECHANICALLY CHECKED PROOF OF IEEE COMPLIANCE OF THE FLOATING POINT MULTIPLICATION, DIVISION AND SQUARE ROOT ALGORITHMS OF THE AMD-K7<sup>TM</sup> PROCESSOR

DAVID M. RUSSINOFF

## *Abstract*

We describe a mechanically verified proof of correctness of the floating point multiplication, division, and square root instructions of the AMD-K7 microprocessor. The instructions are implemented in hardware and represented here by register-transfer level specifications, the primitives of which are logical operations on bit vectors. On the other hand, the statements of correctness, derived from IEEE Standard 754, are arithmetic in nature and considerably more abstract. Therefore, we begin by developing a theory of bit vectors and their role in floating point representations and rounding. We then present the hardware model and a rigorous proof of its correctness. All of our definitions, lemmas and theorems have been formally encoded in the ACL2 logic, and every step in the proof has been mechanically checked with the ACL2 prover.

## 1. *Introduction*

One of the challenges of formal hardware verification is the “semantic gap” between abstract behavioral specifications and concrete hardware models. Dealing effectively with this problem requires a formalism that is flexible enough to represent concepts at different levels of abstraction. In particular, specifications of floating point operations are most naturally expressed in numerical terms, while their hardware implementations are commonly modeled at the level of registers and bit vectors.

Conventional mathematical analysis may be usefully applied to numerical algorithms, but generally fails to provide any assurance regarding the correctness of hardware implementations. On the other hand, automatic finite-state techniques, which have been used to verify low-level specifications of arithmetic circuits [3, 4], lack the expressive power to represent high-level mathematical properties. General-purpose theorem provers offer an important alternative to finite-state tools, as they provide a framework for formal numerical analysis, as well as mechanical support for checking the properties of detailed low-level models.

In our previous work [8] and that of Moore *et al.* [6] on the AMD-K5 floating point unit, the ACL2 theorem prover [2] was used to support the verification of the IEEE compliance [5] of the AMD-K5 floating point division and square root operations. The implementation of these instructions was based on microcode that accessed existing hardware for addition, subtraction, multiplication, and rounding. It was appropriate, therefore, to model the instructions in a language in which the primitive operations included the computation

---

Received 28 January 1998, revised 29 September 1998; published 23 December 1998.

1991 Mathematics Subject Classification 68U30, 65Y99

© 1998, David M. Russinoff

of rounded products and sums, which were assumed to be implemented correctly. Consequently, the analysis was conveniently limited to the familiar realm of floating point numbers and rational arithmetic.

In contrast, the division and square root instructions of the AMD-K7 microprocessor, which were recently designed at AMD by Stuart Oberman [7], are implemented directly in hardware. In order to gain confidence in their correctness, it is desirable to model these instructions at the register-transfer level, where the basic operations are logical functions of bit vectors. Verification then requires bridging the gap between these low-level data and operations and the abstract mathematical objects and functions that they represent.

The subject of this paper is a mechanically verified proof of correctness of the AMD-K7 floating point multiplication, division and square root instructions. The proof is based on a formal description of the hardware, derived from an executable model that was written in C and used for preliminary testing. The instructions are defined in terms of bitwise logical operations and integer addition and multiplication, which are treated as primitives.

The statements of correctness are based on IEEE standard 754 [5], which stipulates that each operation

... shall be performed as if it first produced an intermediate result correct to infinite precision and with unbounded range, and then rounded that result according to one of the [supported] modes ....

Thus, if  $\text{rnd}(x, rc, pc)$  denotes the result of rounding a number  $x$  according to a specified rounding mode  $rc$  and degree of precision  $pc$ , and  $u$  is the value computed for the product of floating point numbers  $a$  and  $b$  in the context of  $rc$  and  $pc$ , then

$$u = \text{rnd}(a \cdot b, rc, pc). \tag{1}$$

Similarly, if  $v$  and  $w$  are the values computed for the quotient of  $a$  and  $b$ , and the square root of  $b$ , respectively, then

$$v = \text{rnd}(a/b, rc, pc) \tag{2}$$

and

$$w = \text{rnd}(\sqrt{b}, rc, pc). \tag{3}$$

The decision to use ACL2, however, has influenced our formulation of this last specification. As a subset of Common Lisp [9], ACL2 includes the rational numbers as a data type, but not the reals. Consequently, we are somewhat limited in our formalization. The reader will notice that many of our lemmas are truths about real numbers but are presented here as propositions of rational arithmetic. More critically, since the square root itself is not a rational function, we are unable to formalize Equation 3 directly. Instead, we prove the following rational version: *For any nonnegative rational numbers  $\ell$  and  $h$ , if  $\ell^2 \leq P \leq h^2$ , then*

$$\text{rnd}(\ell, rc, pc) \leq w \leq \text{rnd}(h, rc, pc). \tag{4}$$

As shown in [8], the equivalence of Equations 3 and 4 is a simple consequence of (a) the monotonicity of rounding, and (b) the observation that for fixed  $rc$  and  $pc$ , the function  $\text{rnd}$  is constant in some neighborhood of any given irrational number.

Applied to the design of a device as complex as a floating point divider, mathematical proof provides a level of confidence that cannot be achieved through testing alone. In the present case, initial proof attempts revealed two design flaws that had survived some

80 million test vectors. The value of mechanical verification in this context is also clear: comprehensive analysis of a commercial floating point design is difficult if not impossible without computer assistance; in any case, the level of investment in its correctness requires a more efficient means of assurance than the conventional social process by which mathematical results are usually confirmed. This is not an argument, however, for circumventing the normal review process. The obligation to support a scientific claim cannot be satisfied simply by announcing that its correctness has been affirmed by an arcane automated proof system, the soundness of which itself is open to question. Moreover, the advantages of a coherent, surveyable proof extend beyond the issue of reliability: it is the only means by which a theory or result may be fully understood, applied, generalized, and assimilated into the mathematical domain. Traditional mathematical notation is clearly a better choice of medium for such an exposition than any formal language.

Since machine-assisted proofs have inherent advantages as well as disadvantages with respect to more traditional methods, we endeavor to combine the benefits of both approaches. In the following sections, we present a detailed proof of correctness, based on elementary mathematics and using only standard terminology and notation. In Section 2, we establish a general theory of floating point numbers, which should be reusable in a wide variety of applications. This is an extension of the theory presented in [8], including some additional properties of the rounding functions, but more significantly, a comprehensive treatment of bit vectors and their role in floating point representation. The specific hardware model is presented in Sections 3 and 4, along with precise formulations and detailed proofs of the above Equations 1, 2 and 4.

For the most part, each step in the proof may be readily checked by hand, requiring no special background in either mathematics or computer hardware. The only exception occurs in Section 4.2, where the accuracy of an approximation derived from a set of tables depends on properties of the tables that can only be verified by extensive (although straightforward) computation, involving approximately  $10^5$  table accesses and  $10^6$  arithmetic operations. The results of these calculations are stated without proof in Lemmas 4.1, 4.2, and 4.3.

On the other hand, along with the table calculations, every step in the proof, including every theorem and lemma presented below, has been formally encoded in the ACL2 logic and mechanically checked with the ACL2 prover, in the interest of eliminating the possibility of human error. The input to the prover, culminating in formal versions of our three main theorems, consisted of some 250 definitions and 3000 lemmas, in addition to the relevant definitions and lemmas of the previously developed general theory [8]. For the interested reader, the files containing this input are included in Appendix A, available to subscribers to the journal at: <http://www.lms.ac.uk/jcm/1/lms98001/appendix-a/>.

## 2. Floating point arithmetic

This section is a formalization of the floating point representation of rational numbers and rounding. The sets of rational numbers, positive rationals, integers, positive integers, and natural numbers (nonnegative integers) will be denoted by the symbols  $\mathbb{Q}$ ,  $\mathbb{Q}^+$ ,  $\mathbb{Z}$ ,  $\mathbb{Z}^+$ , and  $\mathbb{N}$ , respectively. If  $m \in \mathbb{Z}$ ,  $n \in \mathbb{Z}^+$ , and  $m = nq + r$ , where  $q \in \mathbb{Z}$ ,  $r \in \mathbb{N}$ , and  $r < n$ , then we shall write  $\text{rem}(m, n) = r$ .

For  $x \in \mathbb{Q}$ ,  $\lfloor x \rfloor$  and  $\lceil x \rceil$  denote the *floor* and *ceiling* of  $x$ , respectively, defined to be the unique integers satisfying  $\lfloor x \rfloor \leq x < \lfloor x \rfloor + 1$  and  $\lceil x \rceil \geq x > \lceil x \rceil - 1$ . We shall assume familiarity with the basic properties of these functions, including the following.

- (1) If  $n \in \mathbb{Z}$ , then  $\lfloor x + n \rfloor = \lfloor x \rfloor + n$ .
- (2) If  $n \in \mathbb{Z}^+$ , then  $\lfloor \lfloor x \rfloor / n \rfloor = \lfloor x / n \rfloor$ .
- (3) If  $m \in \mathbb{Z}$  and  $n \in \mathbb{Z}^+$ , then  $\lfloor -(m + 1) / n \rfloor = -\lfloor m / n \rfloor - 1$ .

2.1. Bit vectors

We shall exploit the natural correspondence between the bit vectors of length  $n$  and the natural numbers in the range  $0 \leq x < 2^n$ , under which the vector  $b_{n-1}b_{n-2} \cdots b_1b_0$ , where each  $b_k \in \{0, 1\}$ , corresponds to  $x = \sum_{k=0}^{n-1} 2^k b_k$ . The  $k^{th}$  bit of  $x$ ,  $x[k] = b_k$ , is formally defined as follows.

**Definition 2.1.** For all  $x, k \in \mathbb{N}$ ,  $x[k] = \text{rem}(\lfloor x / 2^k \rfloor, 2)$ .

We have the following alternate characterization of  $x[k]$ .

**Lemma 2.1.** For all  $x, k \in \mathbb{N}$ ,  $x[k] = \begin{cases} \text{rem}(x, 2) & \text{if } k = 0 \\ \lfloor x / 2 \rfloor [k - 1] & \text{if } k > 0. \end{cases}$

*Proof.* For  $k > 0$ ,  $x[k] = \text{rem}(\lfloor x / 2^k \rfloor, 2) = \text{rem}(\lfloor \lfloor x / 2 \rfloor / 2^{k-1} \rfloor, 2) = \lfloor x / 2 \rfloor [k - 1]$ . □

**Lemma 2.2.** For all  $x, n, k \in \mathbb{N}$ ,

- (a) if  $x < 2^n$ , then  $x[n] = 0$ ;
- (b) if  $k < n$  and  $2^n - 2^k \leq x < 2^n$ , then  $x[k] = 1$ .

*Proof.* (a)  $x[n] = \text{rem}(\lfloor x / 2^n \rfloor, 2) = \text{rem}(0, 2) = 0$ .

(b) Since  $2^{n-k} - 1 \leq x / 2^k < 2^{n-k}$ ,  $\text{rem}(\lfloor x / 2^k \rfloor, 2) = \text{rem}(2^{n-k} - 1, 2) = 1$ . □

**Lemma 2.3.** For all  $x, m, n \in \mathbb{N}$ ,

- (a)  $(x + 2^n)[n] \neq x[n]$ ;
- (b) if  $m > n$ , then  $\text{rem}(x, 2^m)[n] = x[n]$ .

*Proof.* For any  $m \geq n$  and  $q \in \mathbb{N}$ ,

$$(x + 2^m q)[n] = \text{rem}(\lfloor (x + 2^m q) / 2^n \rfloor, 2) = \text{rem}(\lfloor x / 2^n \rfloor + 2^{m-n} q, 2).$$

If  $m = n$ , then  $\text{rem}(\lfloor x / 2^n \rfloor + 2^{m-n}, 2) = \text{rem}(\lfloor x / 2^n \rfloor + 1, 2) \neq \text{rem}(\lfloor x / 2^n \rfloor, 2) = x[n]$ ; if  $m > n$ , then  $2^{m-n} q$  is even and  $\text{rem}(\lfloor x / 2^n \rfloor + 2^{m-n} q, 2) = \text{rem}(\lfloor x / 2^n \rfloor, 2) = x[n]$ . □

The *left* and *right shift* functions (*shl* and *shr*) take three arguments: a bit vector  $x$ , its length  $n$ , and a value  $s \in \{0, 1\}$  to be shifted in.

**Definition 2.2.** Let  $x, n, s \in \mathbb{N}$  with  $x < 2^n$  and  $s < 2$ .

- (a)  $\text{shl}(x, s, n) = \text{rem}(2x + s, 2^n)$ ;
- (b)  $\text{shr}(x, s, n) = \lfloor x / 2 \rfloor + 2^{n-1} s$ .

*Concatenation (cat)* is also a function of three arguments: two bit vectors,  $x$  and  $y$ , and the length  $n$  of  $y$ .

**Definition 2.3.** For all  $x, y, n \in \mathbb{N}$ ,  $\text{cat}(x, y, n) = 2^n x + y$ .

The following function extracts a field of bits.

**Definition 2.4.** For all  $x, i, j \in \mathbb{N}$ ,  $x[i : j] = \lfloor \text{rem}(x, 2^{j+1}) / 2^i \rfloor$ .

**Lemma 2.4.** For all  $x, y, i, j \in \mathbb{N}$ , if  $\text{rem}(x, 2^{i+1}) = \text{rem}(y, 2^{i+1})$ , then  $x[i : j] = y[i : j]$ .

*Proof.*  $x[i : j] = \lfloor \text{rem}(x, 2^{i+1})/2^j \rfloor = \lfloor \text{rem}(y, 2^{i+1})/2^j \rfloor = y[i : j]$ . □

**Lemma 2.5.** For all  $x, i, j, k, \ell \in \mathbb{N}$ ,

- (a) if  $i \geq k$  and  $j \geq k$ , then  $x[i : j] = \lfloor x/2^k \rfloor [i - k : j - k]$ ;
- (b) if  $i \geq j + k$ , then  $x[i : j][k] = x[k + j]$ ;
- (c) if  $i \geq j + k$ , then  $x[i : j][k : \ell] = x[k + j : \ell + j]$ .

*Proof.* (a) Let  $x = 2^{i+1}q + r$ , where  $0 \leq r < 2^{i+1}$ . Then

$$\lfloor x/2^k \rfloor = \lfloor 2^{i-k+1}q + r/2^k \rfloor = 2^{i-k+1}q + \lfloor r/2^k \rfloor;$$

hence

$$\text{rem}(\lfloor x/2^k \rfloor, 2^{i-k+1}) = \lfloor r/2^k \rfloor$$

and

$$\lfloor x/2^k \rfloor [i - k : j - k] = \lfloor \lfloor r/2^k \rfloor / 2^{j-k} \rfloor = \lfloor r/2^j \rfloor = \lfloor \text{rem}(x, 2^{i+1})/2^j \rfloor = x[i : j].$$

(b) Using Lemma 2.3,

$$\begin{aligned} x[i : j][k] &= \text{rem}(\lfloor \text{rem}(x, 2^{i+1})/2^j \rfloor / 2^k, 2) = \text{rem}(\lfloor \text{rem}(x, 2^{i+1})/2^{k+j} \rfloor, 2) \\ &= \text{rem}(x, 2^{i+1})[k + j] = x[k + j]. \end{aligned}$$

(c) Using (a),

$$\begin{aligned} x[i : j][k : \ell] &= \lfloor x/2^j \rfloor [i - j : 0][k : \ell] = \text{rem}(\lfloor x/2^j \rfloor, 2^{i-j+1})[k : \ell] \\ &= \lfloor \text{rem}(\text{rem}(\lfloor x/2^j \rfloor, 2^{i-j+1}), 2^{k+1})/2^\ell \rfloor = \lfloor \text{rem}(\lfloor x/2^j \rfloor, 2^{k+1})/2^\ell \rfloor \\ &= \lfloor x/2^j \rfloor [k : \ell] = x[k + j : \ell + j]. \end{aligned}$$

□

We have two unary operations on bit vectors, *complement (comp)* and *decrement (dec)*.

**Definition 2.5.** For all  $x, n \in \mathbb{N}$ , if  $x < 2^n$ , then

- (a)  $\text{comp1}(x, n) = 2^n - x - 1$ ;
- (b)  $\text{dec1}(x, n) = \text{rem}(2^n + x - 1, 2^n)$ .

We have three binary logical operations, *and*, *or*, and *exclusive or*.

**Definition 2.6.** For all  $x, y \in \mathbb{N}$ ,

- (a)  $x \ \& \ y = \begin{cases} 0 & \text{if } x = 0 \\ 2(\lfloor x/2 \rfloor \ \& \ \lfloor y/2 \rfloor) + 1 & \text{if } x \text{ and } y \text{ are both odd} \\ 2(\lfloor x/2 \rfloor \ \& \ \lfloor y/2 \rfloor) & \text{otherwise.} \end{cases}$
- (b)  $x \ | \ y = \begin{cases} y & \text{if } x = 0 \\ 2(\lfloor x/2 \rfloor \ | \ \lfloor y/2 \rfloor) & \text{if } x \text{ and } y \text{ are both even} \\ 2(\lfloor x/2 \rfloor \ | \ \lfloor y/2 \rfloor) + 1 & \text{otherwise.} \end{cases}$
- (c)  $x \ \wedge \ y = \begin{cases} y & \text{if } x = 0 \\ 2(\lfloor x/2 \rfloor \ \wedge \ \lfloor y/2 \rfloor) & \text{if } \text{rem}(x, 2) = \text{rem}(y, 2) \\ 2(\lfloor x/2 \rfloor \ \wedge \ \lfloor y/2 \rfloor) + 1 & \text{otherwise.} \end{cases}$

The remainder of this subsection is a collection of properties of the binary logical operations.

**Lemma 2.6.** For all  $x, y \in \mathbb{N}$ ,

- (a)  $x \& y = 2(\lfloor x/2 \rfloor \& \lfloor y/2 \rfloor) + (\text{rem}(x, 2) \& \text{rem}(y, 2));$
- (b)  $x \mid y = 2(\lfloor x/2 \rfloor \mid \lfloor y/2 \rfloor) + (\text{rem}(x, 2) \mid \text{rem}(y, 2)).$

*Proof.* The equivalences are easily checked for all possible values of  $\text{rem}(x, 2)$  and  $\text{rem}(y, 2)$ . □

**Lemma 2.7.** For all  $x, y, z \in \mathbb{N}$ ,

- (a)  $x \& 0 = 0;$
- (b)  $x \mid 0 = x;$
- (c)  $x \& y = y \& x;$
- (d)  $x \mid y = y \mid x;$
- (e)  $(x \& y) \& z = x \& (y \& z);$
- (f)  $(x \mid y) \mid z = x \mid (y \mid z);$
- (g)  $x \mid (y \& z) = (x \mid y) \& (x \mid z);$
- (h)  $x \& (y \mid z) = (x \& y) \mid (x \& z).$

*Proof.* First note that Lemma 2.6 implies

$$\lfloor (x \& y)/2 \rfloor = \lfloor x/2 \rfloor \& \lfloor y/2 \rfloor \text{ and } \text{rem}(x \& y, 2) = \text{rem}(x, 2) \& \text{rem}(y, 2)$$

and

$$\lfloor (x \mid y)/2 \rfloor = \lfloor x/2 \rfloor \mid \lfloor y/2 \rfloor \text{ and } \text{rem}(x \mid y, 2) = \text{rem}(x, 2) \mid \text{rem}(y, 2).$$

We shall prove (h); the other proofs are similar.

It is easily verified that the statement holds for arguments in  $\{0, 1\}$ . Thus,

$$\begin{aligned} \text{rem}(x \& (y \mid z), 2) &= \text{rem}(x, 2) \& \text{rem}(y \mid z, 2) \\ &= \text{rem}(x, 2) \& (\text{rem}(y, 2) \mid \text{rem}(z, 2)) \\ &= (\text{rem}(x, 2) \& \text{rem}(y, 2)) \mid (\text{rem}(x, 2) \& \text{rem}(z, 2)) \\ &= \text{rem}(x \& y, 2) \mid (\text{rem}(x \& z, 2)) \\ &= \text{rem}((x \& y) \mid (x \& z), 2). \end{aligned}$$

Now, by inductive hypothesis,

$$\begin{aligned} \lfloor (x \& (y \mid z))/2 \rfloor &= \lfloor x/2 \rfloor \& \lfloor (y \mid z)/2 \rfloor \\ &= \lfloor x/2 \rfloor \& (\lfloor y/2 \rfloor \mid \lfloor z/2 \rfloor) \\ &= (\lfloor (x \& y)/2 \rfloor) \mid (\lfloor (x \& z)/2 \rfloor) \\ &= (\lfloor x/2 \rfloor \& \lfloor y/2 \rfloor) \mid (\lfloor x/2 \rfloor \& \lfloor z/2 \rfloor) \\ &= \lfloor ((x \& y) \mid (x \& z))/2 \rfloor. \end{aligned}$$

Therefore,

$$\begin{aligned} x \& (y \mid z) &= \lfloor (x \& (y \mid z))/2 \rfloor + \text{rem}(x \& (y \mid z), 2) \\ &= \lfloor ((x \& y) \mid (x \& z))/2 \rfloor + \text{rem}((x \& y) \mid (x \& z), 2) \\ &= (x \& y) \mid (x \& z). \end{aligned}$$

□

**Lemma 2.8.** *Let  $x, y, n \in \mathbb{N}$ .*

- (a) *if  $x < 2^n$  and  $y < 2^n$ , then  $x \mid y < 2^n$ ;*
- (b) *if  $y < 2^n$ , then  $(2^n x) \mid y = 2^n x + y$ ;*
- (c)  *$(2^n x) \mid (2^n y) = 2^n(x \mid y)$ ;*
- (d)  *$\text{rem}(x \mid y, 2^n) = \text{rem}(x, 2^n) \mid \text{rem}(y, 2^n)$ .*

*Proof.* (a) For  $n > 0$ ,  $\lfloor x/2 \rfloor < 2^{n-1}$  and  $\lfloor y/2 \rfloor < 2^{n-1}$ , which implies  $\lfloor x/2 \rfloor \mid \lfloor y/2 \rfloor < 2^{n-1}$ ; hence

$$x \mid y \leq 2(\lfloor x/2 \rfloor \mid \lfloor y/2 \rfloor) + 1 \leq 2(2^{n-1} - 1) + 1 < 2^n.$$

(b) For  $n > 0$ , since  $\lfloor y/2 \rfloor < 2^{n-1}$ ,

$$\begin{aligned} (2^n x) \mid y &= 2(\lfloor 2^n x/2 \rfloor \mid \lfloor y/2 \rfloor) + \text{rem}(2^n x, 2) \mid \text{rem}(y, 2) \\ &= 2(2^{n-1} x \mid \lfloor y/2 \rfloor) + 0 \mid \text{rem}(y, 2) \\ &= 2(2^{n-1} x + \lfloor y/2 \rfloor) + \text{rem}(y, 2) \\ &= 2^n x + 2\lfloor y/2 \rfloor + \text{rem}(y, 2) \\ &= 2^n x + y. \end{aligned}$$

(c) For  $n > 0$ ,

$$\begin{aligned} (2^n x) \mid (2^n y) &= 2(\lfloor 2^n x/2 \rfloor \mid \lfloor 2^n y/2 \rfloor) + \text{rem}(2^n x, 2) \mid \text{rem}(2^n y, 2) \\ &= 2(2^{n-1} x \mid 2^{n-1} y) + 0 \mid 0 = 2(2^{n-1}(x \mid y)) + 0 \\ &= 2^n(x \mid y). \end{aligned}$$

(d) Let  $x = 2^n q_1 + r_1$  and  $y = 2^n q_2 + r_2$ , where  $0 \leq r_1 < 2^n$  and  $0 \leq r_2 < 2^n$ . Then

$$\begin{aligned} x \mid y &= (2^n q_1 + r_1) \mid (2^n q_2 + r_2) = (2^n q_1 \mid r_1) \mid (2^n q_2 \mid r_2) \\ &= (2^n q_1 \mid 2^n q_2) \mid (r_1 \mid r_2) = (2^n(q_1 \mid q_2)) \mid (r_1 \mid r_2) \\ &= 2^n(q_1 \mid q_2) + (r_1 \mid r_2). \end{aligned}$$

But  $r_1 \mid r_2 < 2^n$ ; hence  $\text{rem}(x \mid y, 2^n) = r_1 \mid r_2 = \text{rem}(x, 2^n) \mid \text{rem}(y, 2^n)$ . □

**Lemma 2.9.** *Let  $x, y, n \in \mathbb{N}$ .*

- (a)  $x \ \& \ y \leq x$ ;
- (b)  $2^n x \ \& \ y = 2^n(x \ \& \ \lfloor y/2^n \rfloor)$ ;
- (c)  $\text{rem}(x \ \& \ y, 2^n) = \text{rem}(x, 2^n) \ \& \ y$ ;
- (d) *if  $x < 2^n$ , then  $x \ \& \ y = x \ \& \ \text{rem}(y, 2^n)$ .*

*Proof.* (a) If  $x = 0$ , then  $x \ \& \ y = 0 \leq x$ , and for  $x > 0$ ,

$$\begin{aligned} x \ \& \ y &= 2(\lfloor x/2 \rfloor \ \& \ \lfloor y/2 \rfloor) + (\text{rem}(x, 2) \ \& \ \text{rem}(y, 2)) \leq 2\lfloor x/2 \rfloor + \text{rem}(x, 2) \\ &= x. \end{aligned}$$

(b) For  $n > 0$ ,

$$\begin{aligned} 2^n x \ \& \ y &= 2(\lfloor 2^n x/2 \rfloor \ \& \ \lfloor y/2 \rfloor) + \text{rem}(2^n x, 2) \ \& \ \text{rem}(y, 2) \\ &= 2(2^{n-1} x \ \& \ \lfloor y/2 \rfloor) + 0 \ \& \ \text{rem}(y, 2) \\ &= 2(2^{n-1}(x \ \& \ \lfloor \lfloor y/2 \rfloor / 2^{n-1} \rfloor)) + 0 \\ &= 2^n(x \ \& \ \lfloor y/2^n \rfloor). \end{aligned}$$

(c) Let  $x = 2^n q + r$ ,  $0 \leq r < 2^n$ . Then  $0 \leq r$  &  $y \leq r < 2^n$  and

$$\begin{aligned} x \& y &= (2^n q + r) \& y = (2^n q \mid r) \& y \\ &= (2^n q \& y) \mid (r \& y) = (2^n (q \& \lfloor y/2^n \rfloor) \mid (r \& y)) \\ &= (2^n (q \& \lfloor y/2^n \rfloor) + (r \& y)). \end{aligned}$$

Therefore,  $\text{rem}(x \& y, 2^n) = r \& y = \text{rem}(x, 2^n) \& y$ .

(d) Since  $x \& y \leq x < 2^n$ ,  $x \& y = \text{rem}(x \& y, 2^n) = x \& \text{rem}(y, 2^n)$ . □

**Lemma 2.10.** Let  $x, y, n \in \mathbb{N}$ .

- (a)  $(x \& y)[n] = x[n] \& y[n]$ ;
- (b)  $(x \mid y)[n] = x[n] \mid y[n]$ .

*Proof.* The proofs are similar; we present the proof of (a), which proceeds by induction. For  $n = 0$ ,

$$(x \& y)[0] = \text{rem}(x \& y, 2) = \text{rem}(x, 2) \& \text{rem}(y, 2) = x[0] \& y[0];$$

for  $n > 0$ ,

$$\begin{aligned} (x \& y)[n] &= \lfloor (x \& y)/2 \rfloor [n - 1] = (\lfloor x/2 \rfloor \& \lfloor y/2 \rfloor)[n - 1] \\ &= \lfloor x/2 \rfloor [n - 1] \& \lfloor y/2 \rfloor [n - 1] = x[n] \& y[n]. \end{aligned}$$

□

**Lemma 2.11.** Let  $x, n, k \in \mathbb{N}$ ,  $k < n$ .

- (a)  $x \& 2^k = 2^k x[k]$ ;
- (b)  $x \mid 2^k = x + 2^k(1 - x[k])$ ;
- (c)  $x \& (2^n - 2^k) = 2^k(x[n - 1 : k])$ .

*Proof.* (a) In the case  $k = 0$ , we have

$$x \& 1 = 2(\lfloor x/2 \rfloor \& 0) + \text{rem}(x, 2) = \text{rem}(x, 2) = x[0],$$

and for  $k > 0$ , by Lemma 2.1,

$$x \& 2^k = 2(\lfloor x/2 \rfloor \& 2^{k-1}) = 2(2^{k-1} \lfloor x/2 \rfloor [k - 1]) = 2^k x[k].$$

(b) For  $k = 0$ , we have

$$x \mid 1 = 2(\lfloor x/2 \rfloor \mid 0) + 1 = 2\lfloor x/2 \rfloor + 1 = x + 1 - \text{rem}(x, 2) = x + 1 - x[0],$$

and for  $k > 0$ ,

$$\begin{aligned} x \mid 2^k &= 2\{\lfloor x/2 \rfloor \mid 2^{k-1}\} + \text{rem}(x, 2) \\ &= 2\left\{\lfloor x/2 \rfloor + 2^{k-1}(1 - \lfloor x/2 \rfloor [k - 1])\right\} + \text{rem}(x, 2) \\ &= 2\lfloor x/2 \rfloor + \text{rem}(x, 2) + 2^k(1 - \lfloor x/2 \rfloor [k - 1]) \\ &= x + 2^k(1 - x[k]). \end{aligned}$$

(c) It suffices to prove the identity under the assumption  $x < 2^n$ , because then, by Lemmas 2.9 and 2.4, we have for arbitrary  $x$ :

$$x \& (2^n - 2^k) = \text{rem}(x, 2^n) \& (2^n - 2^k) = \text{rem}(x, 2^n)[n : k] = x[n : k].$$



For  $k = 0$ , we show by induction that  $x \& (2^n - 1) = x$ . The case  $n = 0$  is trivial, and for  $n > 0$ , since  $\lfloor (2^n - 1)/2 \rfloor = 2^{n-1} - 1$ , we have

$$\begin{aligned} x \& (2^n - 1) &= 2(\lfloor x/2 \rfloor \& (2^{n-1} - 1)) + \text{rem}(x, 2) \\ &= 2\lfloor x/2 \rfloor + \text{rem}(x, 2) = x. \end{aligned}$$

Now, for  $k > 0$ ,

$$\begin{aligned} x \& (2^n - 2^k) &= 2(\lfloor x/2 \rfloor \& (2^{n-1} - 2^{k-1})) = 2 \cdot 2^{k-1} \lfloor x/2 \rfloor [n - 2 : k - 1] \\ &= 2^k \lfloor \text{rem}(\lfloor x/2 \rfloor, 2^{n-1})/2^{k-1} \rfloor = 2^k \lfloor \lfloor x/2 \rfloor / 2^{k-1} \rfloor \\ &= 2^k \lfloor x/2^k \rfloor = 2^k (x[n - 1 : k]). \end{aligned}$$

□

**Lemma 2.12.** Let  $n, k, \ell \in \mathbb{N}$ ,  $\ell \leq k < n$ . Then

$$(2^n - 2^\ell - 1) \& (2^n - 2^k) = \begin{cases} 2^n - 2^{k+1} & \text{if } \ell = k \\ 2^n - 2^k & \text{if } \ell < k. \end{cases}$$

*Proof.* Applying Lemma 2.11 (c), we have

$$\begin{aligned} (2^n - 2^\ell - 1) \& (2^n - 2^k) &= 2^k (2^n - 2^\ell - 1)[n - 1 : k] = 2^k \lfloor (2^n - 2^\ell - 1)/2^k \rfloor \\ &= 2^k (2^{n-k} + \lfloor -(2^\ell + 1)/2^k \rfloor) \\ &= 2^n - 2^k (\lfloor 2^{\ell-k} \rfloor + 1). \end{aligned}$$

□

## 2.2. Floating point representations

Floating point representation is based on the observation that every nonzero rational number  $x$  admits a unique factorization,

$$x = \text{sgn}(x) \text{sig}(x) 2^{\text{expo}(x)},$$

where  $\text{sgn}(x) \in \{1, -1\}$  (the *sign* of  $x$ ),  $1 \leq \text{sig}(x) < 2$  (the *significand* of  $x$ ), and  $\text{expo}(x) \in \mathbb{Z}$  (the *exponent* of  $x$ ).

**Definition 2.7.** Let  $x \in \mathbb{Q}$ . If  $x \neq 0$ , then

- (a)  $\text{sgn}(x) = x/|x|$ ;
- (b)  $\text{expo}(x)$  is the unique integer that satisfies  $2^{\text{expo}(x)} \leq |x| < 2^{\text{expo}(x)+1}$ ;
- (c)  $\text{sig}(x) = |x| 2^{-\text{expo}(x)}$ .

A floating point representation of  $x$  consists of three bit vectors, corresponding to  $\text{sgn}(x)$ ,  $\text{sig}(x)$ , and  $\text{expo}(x)$ . A format is defined by the number of bits allocated to  $\text{sig}(x)$  and  $\text{expo}(x)$ .

**Definition 2.8.** Let  $\phi = (\mu, \epsilon) \in \mathbb{Z}^+ \times \mathbb{Z}^+$ . Then  $\phi$  is a floating point format. A  $\phi$ -encoding is a triple  $(s, m, e) \in \mathbb{N} \times \mathbb{N} \times \mathbb{N}$  such that  $s < 2$ ,  $m < 2^\mu$ , and  $e < 2^\epsilon$ . If  $z = (s, m, e)$ , then  $s = \text{get-sign}(z)$ ,  $m = \text{get-man}(z)$ , and  $e = \text{get-expo}(z)$ . If  $m \geq 2^{\mu-1}$ , then  $z$  is a *normal*  $\phi$ -encoding.

The formats that are supported by the AMD-K7 floating point operations include (24, 7), (53, 10), and (64, 15), which correspond to *single*, *double*, and *extended* precision as specified by IEEE, as well as a larger format, (68, 18). In addition, in order to allow for the rounding error incurred by our iterative division and square root algorithms, which are required to produce results that are correctly rounded to 68 bits, the multiplier must support a somewhat more precise internal format. One of the objectives of our analysis is to determine the minimum required size of this format, and hence the minimum width of the multiplier. Thus, we introduce an integer parameter  $M$ , which represents the multiplier width and determines the internal format  $(M, 18)$ . We assume that  $M \geq 75$ , for as we shall see in Section 4, our proofs of correctness for division and square root will depend on this constraint.

In our formulation of the algorithms, the floating point formats are encoded as symbols.

**Definition 2.9.** A *precision control specifier* is any of the symbols

$$\text{PC-32, PC-64, PC-80, PC-87, and PC-*,}$$

which correspond to the floating point formats

$$(24, 7), (53, 10), (64, 15), (68, 18), \text{ and } (M, 18),$$

respectively. The first four of these symbols are called *external* precision control specifiers. If  $\pi$  is any precision control specifier and  $\phi = (\mu, \epsilon)$  is the corresponding format, then  $\text{mbits}(\pi) = \mu$ .

The number  $x$  represented by a normal  $(\mu, \epsilon)$ -encoding  $(s, m, e)$  is given by  $\text{sgn}(x) = (-1)^s$ ,  $\text{sig}(x) = 2^{\mu-1}m$ , and  $\text{expo}(x) = e - (2^{\epsilon-1} - 1)$ . Thus, the exponent is biased in order to provide for an exponent range  $1 - 2^{\epsilon-1} \leq \text{expo}(x) \leq 2^{\epsilon-1}$ .

**Definition 2.10.** Let  $z = (s, m, e)$  be a  $\phi$ -encoding, where  $\phi = (\mu, \epsilon)$  is a floating point format. Then  $\text{decode}(z, \phi) = (-1)^s \cdot m \cdot 2^{e-2^{\epsilon-1}-\mu+2}$ . In the case  $\phi = (M, 18)$ , we shall designate  $x$  simply as an *encoding*, and  $\text{decode}(x, (M, 18))$  will be denoted as  $\hat{x}$ .

Our characterization of the rational numbers that are represented by normal encodings is based on the following definition.

**Definition 2.11.** Let  $x \in \mathbb{Q}$  and  $n \in \mathbb{Z}^+$ . Then  $x$  is *n-exact* iff  $\text{sig}(x)2^{n-1} \in \mathbb{Z}$ .

The following basic property of *n-exact* numbers is proved in [8].

**Lemma 2.13.** If  $x \in \mathbb{Q}^+$ ,  $n \in \mathbb{Z}^+$ , and  $x$  is *n-exact*, then the least *n-exact* number that is greater than  $x$  is  $x + 2^{\text{expo}(x)+1-n}$ .

We shall also require this trivial characterization of *n-exact* bit vectors.

**Lemma 2.14.** Let  $x, n, k \in \mathbb{Z}^+$ ,  $2^{n-1} \leq x < 2^n$  and  $k < n$ . The following are equivalent.

- (a)  $2^k$  divides  $x$ ;
- (b)  $x$  is  $(n - k)$ -exact;
- (c)  $x[n - 1 : k] = x/2^k$ ;
- (d)  $x[k - 1 : 0] = 0$ .

**Definition 2.12.** Let  $x \in \mathbb{Q}$  and let  $\phi = (\mu, \epsilon)$  be a floating point format. Then  $x$  is  $\phi$ -representable iff  $x$  is  $\mu$ -exact and  $-2^{\epsilon-1} + 1 \leq \text{expo}(x) \leq 2^{\epsilon-1}$ . If  $\phi = (M, 18)$ , then we shall say that  $x$  is *representable*.

The inverse of decode is given below.

**Definition 2.13.** Let  $\phi = (\mu, \epsilon)$  be a floating point format and let  $x$  be  $\phi$ -representable,  $x \neq 0$ . Then  $\text{encode}(x, \phi) = (s, m, e)$ , where

- (a) if  $\text{sgn}(x) = 1$ , then  $s = 0$ , and if  $\text{sgn}(x) = -1$ , then  $s = 1$ ;
- (b)  $m = \text{sig}(x)2^{\mu-1}$ ;
- (c)  $e = \text{expo}(x) + 2^{\epsilon-1} - 1$ .

**Lemma 2.15.** Let  $\phi = (\mu, \epsilon)$  be a floating point format, let  $z = (s, m, e)$  be a normal  $\phi$ -encoding, and let  $x = \text{decode}(z, \phi)$ .

- (a)  $\text{sgn}(x) = (-1)^s$ ;
- (b)  $\text{sig}(x) = m/2^{\mu-1}$ ;
- (d)  $x$  is  $\phi$ -representable;
- (c)  $\text{expo}(x) = e - 2^{\epsilon-1} + 1$ ;
- (e)  $\text{encode}(x, \phi) = z$ .

*Proof.* Let  $\phi = (\mu, \epsilon)$ . Then

$$x = (-1)^s m 2^{e-(2^{\epsilon-1}-1)-\mu+1} = (-1)^s (m 2^{1-\mu}) 2^{e-(2^{\epsilon-1}-1)}.$$

But  $2^{\mu-1} \leq m < 2^\mu$  yields  $1 \leq m 2^{1-\mu} < 2$ , which implies (a), (b), and (c). Now (d) follows from the relation  $0 \leq e < 2^\epsilon$ , and (e) from Definition 2.13.  $\square$

### 2.3. Rounding

A *rounding mode* is a function  $\mathcal{M}$  that computes an  $n$ -exact number  $\mathcal{M}(x, n)$  corresponding to an arbitrary rational  $x$  and a degree of precision  $n \in \mathbb{Z}^+$ . We define five rounding modes.

**Definition 2.14.** A *rounding mode* is any of the functions *trunc*, *away*, *near*, *inf*, and *minf*, where, for  $x \in \mathbb{Q}$  and  $n \in \mathbb{Z}^+$ ,

- (a)  $\text{trunc}(x, n) = \text{sgn}(x) \lfloor 2^{n-1} \text{sig}(x) \rfloor 2^{\text{expo}(x)-n+1}$ ;
- (b)  $\text{away}(x, n) = \text{sgn}(x) \lceil 2^{n-1} \text{sig}(x) \rceil 2^{\text{expo}(x)-n+1}$ ;
- (c) if  $z = \lfloor 2^{n-1} \text{sig}(x) \rfloor$  and  $f = 2^{n-1} \text{sig}(x) - z$ , then

$$\text{near}(x, n) = \begin{cases} \text{trunc}(x, n) & \text{if } f < 1/2 \\ \text{away}(x, n) & \text{if } f > 1/2 \\ \text{trunc}(x, n) & \text{if } f = 1/2 \text{ and } z \text{ is even} \\ \text{away}(x, n) & \text{if } f = 1/2 \text{ and } z \text{ is odd;} \end{cases}$$

- (d)  $\text{inf}(x, n) = \begin{cases} \text{away}(x, n) & \text{if } x \geq 0 \\ \text{trunc}(x, n) & \text{if } x < 0; \end{cases}$
- (e)  $\text{minf}(x, n) = \begin{cases} \text{trunc}(x, n) & \text{if } x \geq 0 \\ \text{away}(x, n) & \text{if } x < 0. \end{cases}$

Only four of these modes are supported by the IEEE standard. In our representation of the algorithms, they will be encoded as symbols.

**Definition 2.15.** A *rounding control specifier* is any of the symbols

RC-CHOP, RC-POS, RC-NEG, and RC-NEAR,

which correspond to the rounding modes

trunc, inf, minf, and near,

respectively. Let  $\rho$  be a rounding control specifier corresponding to the rounding mode  $\mathcal{M}$ , let  $\pi$  be a precision control specifier, and let  $x \in \mathbb{Q}$ . Then

$$\text{rnd}(x, \rho, \pi) = \mathcal{M}(x, \text{mbits}(\pi)).$$

Some of the basic properties of the rounding modes, which are proved in [8], are listed in the following eight lemmas.

**Lemma 2.16.** *If  $x \in \mathbb{Q}$ ,  $\mathcal{M}$  is a rounding mode, and  $n \in \mathbb{Z}^+$ , then*

- (a)  $\text{sgn}(\mathcal{M}(x, n)) = \text{sgn}(x)$ ;
- (b) if  $\mathcal{M} \in \{\text{trunc, away, near}\}$ , then  $\mathcal{M}(-x, n) = -\mathcal{M}(x, n)$ .

**Lemma 2.17.** *If  $x, y \in \mathbb{Q}$ ,  $x \leq y$ ,  $\mathcal{M}$  is a rounding mode, and  $n \in \mathbb{Z}^+$ , then*

$$\mathcal{M}(x, n) \leq \mathcal{M}(y, n).$$

**Lemma 2.18.** *If  $x \in \mathbb{Q}$ ,  $\mathcal{M}$  is a rounding mode, and  $n \in \mathbb{Z}^+$ , then*

- (a)  $\mathcal{M}(x, n)$  is  $n$ -exact;
- (b) if  $x$  is  $n$ -exact, then  $x = \mathcal{M}(x, n)$ .

**Lemma 2.19.** *If  $x \in \mathbb{Q}$ ,  $\mathcal{M}$  is a rounding mode other than near,  $m, n \in \mathbb{Z}^+$ , and  $m \leq n$ , then*

$$\mathcal{M}(\mathcal{M}(x, n), m) = \mathcal{M}(x, m).$$

**Lemma 2.20.** *If  $x \in \mathbb{Q}$  and  $n \in \mathbb{Z}^+$ , then*

$$|x| - 2^{\text{expo}(x)-n+1} < |\text{trunc}(x, n)| \leq |x| \leq |\text{away}(x, n)| < |x| + 2^{\text{expo}(x)-n+1}.$$

**Lemma 2.21.** *If  $x \in \mathbb{Q}$  and  $n \in \mathbb{Z}^+$ , then*

- (a)  $\text{expo}(\text{trunc}(x, n)) = \text{expo}(x)$ ;
- (b)  $\text{expo}(\text{away}(x, n)) = \text{expo}(x)$  unless  $|\text{away}(x, n)| = 2^{\text{expo}(x)+1}$ .

**Lemma 2.22.** *If  $x, a \in \mathbb{Q}$ ,  $n \in \mathbb{Z}^+$ , and  $a$  is  $n$ -exact, then*

- (a) if  $a \leq |x|$ , then  $a \leq |\text{trunc}(x, n)|$ ;
- (b) if  $a \geq |x|$ , then  $a \geq |\text{away}(x, n)|$ .

**Lemma 2.23.** *Let  $x, y \in \mathbb{Q}$  and  $n \in \mathbb{Z}^+$ . If  $y$  is  $n$ -exact, then  $|x - y| \geq |x - \text{near}(x, n)|$ .*

We shall require a number of properties in addition to the above. The next lemma provides an implementation of truncation of bit vectors.

**Lemma 2.24.** *Let  $x, m, n, k \in \mathbb{N}$ . If  $0 < k < n \leq m$  and  $2^{n-1} \leq x < 2^n$ , then*

$$\text{trunc}(x, k) = x \ \& \ (2^m - 2^{n-k}).$$

*Proof.* By Lemma 2.11,

$$\begin{aligned} \text{trunc}(x, k) &= \lfloor 2^{k-1-\text{expo}(x)} x \rfloor 2^{\text{expo}(x)+1-k} = \lfloor x/2^{n-k} \rfloor 2^{n-k} \\ &= 2^{n-k}(x[n-1 : n-k]) = x \ \& \ (2^n - 2^{n-k}). \end{aligned}$$

But by Lemma 2.9,

$$x \ \& \ (2^m - 2^{n-k}) = x \ \& \ \text{rem}(2^m - 2^{n-k}, 2^n) = x \ \& \ (2^n - 2^{n-k}).$$

□

Lemma 2.24 is also the basis for our implementations of the other rounding modes, which therefore must be characterized in terms of truncation.

**Lemma 2.25.** *Let  $x \in \mathbb{Q}^+$ ,  $m \in \mathbb{Z}^+$ , and  $n \in \mathbb{Z}^+$ . If  $x$  is  $m$ -exact and  $m \geq n$ , then*

$$\text{away}(x, n) = \text{trunc}(x + 2^{\text{expo}(x)+1}(2^{-n} - 2^{-m}), n).$$

*Proof.* Let  $a = \text{trunc}(x + 2^{\text{expo}(x)+1}(2^{-n} - 2^{-m}), n)$ . Since

$$a < x + 2^{\text{expo}(x)+1-n} \leq \text{away}(x, n) + 2^{\text{expo}(\text{away}(x,n))+1-n},$$

$a \leq \text{away}(x, n)$  by Lemma 2.13.

If  $x$  is  $n$ -exact, then  $a \geq \text{trunc}(x, n) = x = \text{away}(x, n)$ , and hence  $a = \text{away}(x, n)$ . Thus, we may assume  $x$  is not  $n$ -exact. But then since  $x > \text{trunc}(x, n)$  and  $x$  is  $m$ -exact,

$$x \geq \text{trunc}(x, n) + 2^{\text{expo}(x)+1-m}$$

and hence

$$x + 2^{\text{expo}(x)+1}(2^{-n} - 2^{-m}) \geq \text{trunc}(x, n) + 2^{\text{expo}(x)+1-n} = \text{away}(x, n),$$

which implies  $a \geq \text{away}(x, n)$ . □

The remainder of this section addresses the properties of *near rounding*, concluding with its characterization as a truncated sum.

**Lemma 2.26.** *If  $x \in \mathbb{Q}$  and  $n \in \mathbb{Z}^+$ , then  $|x - \text{near}(x, n)| \leq 2^{\text{expo}(x)-n}$ .*

*Proof.* By Lemma 2.16, we may assume  $x > 0$ . Let  $a = \text{trunc}(x, n) + 2^{\text{expo}(x)+1-n}$ . By Lemmas 2.18 and 2.23, if the statement fails, then

$$\text{trunc}(x, n) < x - 2^{\text{expo}(x)-n} < x + 2^{\text{expo}(x)-n} < \text{away}(x, n);$$

hence  $a < \text{away}(x, n)$ . Then by Lemmas 2.13 and 2.22(a), we have  $a < x$ , contradicting Lemma 2.22(b). □

**Lemma 2.27.** *Let  $x \in \mathbb{Q}$  and  $n \in \mathbb{Z}^+$ . If  $x$  is  $(n + 1)$ -exact but not  $n$ -exact, then*

- (a)  $\text{trunc}(x, n) = x - \text{sgn}(x)2^{\text{expo}(x)-n}$ ;
- (b)  $\text{away}(x, n) = x + \text{sgn}(x)2^{\text{expo}(x)-n}$ .

*Proof.* Again we may assume  $x > 0$ . Let  $a = x - 2^{\text{expo}(x)-n}$  and  $b = x + 2^{\text{expo}(x)-n}$ . Since  $x > 2^{\text{expo}(x)}$ ,  $x \geq 2^{\text{expo}(x)} + 2^{\text{expo}(x)+1-n}$  by Lemma 2.13; hence  $a \geq 2^{\text{expo}(x)}$  and  $\text{expo}(a) = \text{expo}(x)$ .

By hypothesis,  $x2^{n-\text{expo}(x)}$  is odd. Let  $x2^{n-\text{expo}(x)} = 2k + 1$ . Then

$$a2^{n-1-\text{expo}(a)} = (x - 2^{\text{expo}(x)-n})2^{n-1-\text{expo}(x)} = (2k + 1)/2 - 1/2 = k \in \mathbb{Z}.$$

Thus,  $a$  is  $n$ -exact, and by Lemma 2.13, so is  $a + 2^{\text{expo}(a)+1-n} = b$ . Now by Lemma 2.22,  $a \leq \text{trunc}(x, n)$ , but if  $a < \text{trunc}(x, n)$ , then Lemma 2.13 would imply  $b \leq \text{trunc}(x, n)$ , contradicting  $x < b$ . This establishes (a), and the proof of (b) is similar.  $\square$

**Lemma 2.28.** Let  $x, a \in \mathbb{Q}^+$ , and  $n \in \mathbb{Z}^+$ . If  $a$  is  $n$ -exact, then

- (a) if  $x > a + 2^{\text{expo}(a)-n}$ , then  $\text{near}(x, n) \geq a + 2^{\text{expo}(a)+1-n}$ ;
- (b) if  $x < a + 2^{\text{expo}(a)-n}$ , then  $\text{near}(x, n) \leq a$ ;
- (c) if  $x > a - 2^{\text{expo}(x)-n}$ , then  $\text{near}(x, n) \geq a$ .

*Proof.* (a) Let  $b = a + 2^{\text{expo}(a)+1-n}$ . If  $\text{near}(x, n) < b$ , then Lemma 2.13 yields  $\text{near}(x, n) \leq a$ ; hence  $|\text{near}(x, n) - x| > |\text{near}(x, n) - b|$ , contradicting Lemma 2.23.

(b) If  $\text{near}(x, n) > a$ , then  $\text{near}(x, n) \geq b$ , and a contradiction may be derived as in (a).

(c) By Lemma 2.17, we may assume  $x < a$ . Let  $c = a - 2^{\text{expo}(x)+1-n}$ . Then  $c < x < a$ . Since  $a > x \geq 2^{\text{expo}(x)}$ ,  $a \geq 2^{\text{expo}(x)} + 2^{\text{expo}(x)+1-n}$ , and hence  $x > c \geq 2^{\text{expo}(x)}$ , which implies  $\text{expo}(c) = \text{expo}(x)$ . But  $\text{expo}(c) \leq \text{expo}(a)$  and therefore

$$c2^{n-1-\text{expo}(c)} = a2^{n-1-\text{expo}(c)} - 1 \in \mathbb{Z},$$

i.e.,  $c$  is  $n$ -exact. Now since  $x > a - 2^{\text{expo}(x)-n} = c + 2^{\text{expo}(c)-n}$ , (a) implies  $\text{near}(x, n) \geq c + 2^{\text{expo}(c)+1-n} = a$ .  $\square$

**Lemma 2.29.** Let  $n \in \mathbb{Z}$ ,  $n > 1$ , and  $x \in \mathbb{Q}$ . If  $x$  is  $(n + 1)$ -exact but not  $n$ -exact, then  $\text{near}(x, n)$  is  $(n - 1)$ -exact.

*Proof.* Again we may assume  $x > 0$ . Let  $z = \lfloor 2^{n-1}\text{sig}(x) \rfloor$  and  $f = 2^{n-1}\text{sig}(x) - z$ . Since  $2^{n-1}\text{sig}(x) \notin \mathbb{Z}$ ,  $0 < f < 1$ . But  $2^n\text{sig}(x) = 2z + 2f \in \mathbb{Z}$ ; hence  $2f \in \mathbb{Z}$  and  $f = \frac{1}{2}$ .

If  $z$  is even, then

$$\text{near}(x, n) = \text{trunc}(x, n) = z2^{\text{expo}(x)+1-n}$$

and by Lemma 2.21,

$$2^{n-2-\text{expo}(\text{near}(x,n))}\text{near}(x, n) = 2^{n-2-\text{expo}(x)}z2^{\text{expo}(x)+1-n} = z/2 \in \mathbb{Z}.$$

If  $z$  is odd, then

$$\text{near}(x, n) = \text{away}(x, n) = (z + 1)2^{\text{expo}(x)+1-n}.$$

We may assume  $\text{away}(x, n) \neq 2^{\text{expo}(x)+1}$ ; hence by Lemma 2.21,

$$2^{n-2-\text{expo}(\text{near}(x,n))}\text{near}(x, n) = 2^{n-2-\text{expo}(x)}(z + 1)2^{\text{expo}(x)+1-n} = (z + 1)/2 \in \mathbb{Z}.$$

$\square$

**Lemma 2.30.** Let  $n \in \mathbb{Z}$ ,  $n > 1$ , and  $x \in \mathbb{Q}^+$ . If  $x + 2^{\text{expo}(x)-n} \geq 2^{\text{expo}(x)+1}$ , then

$$\text{near}(x, n) = 2^{\text{expo}(x)+1} = \text{trunc}(x + 2^{\text{expo}(x)-n}, n).$$

*Proof.* Suppose  $\text{near}(x, n) \neq 2^{\text{expo}(x)+1}$ . Then Lemma 2.21 implies  $\text{near}(x, n) < 2^{\text{expo}(x)+1}$  and by Lemmas 2.13 and 2.26,

$$\begin{aligned} 2^{\text{expo}(x)+1} &\geq \text{near}(x, n) + 2^{\text{expo}(x)+1-n} \geq x - 2^{\text{expo}(x)-n} + 2^{\text{expo}(x)+1-n} \\ &= x + 2^{\text{expo}(x)-n} \geq 2^{\text{expo}(x)+1}. \end{aligned}$$

It follows that  $x = 2^{\text{expo}(x)+1} - 2^{\text{expo}(x)-n}$  is  $(n+1)$ -exact but not  $n$ -exact, while  $\text{near}(x, n) = 2^{\text{expo}(x)+1} - 2^{\text{expo}(x)+1-n}$  is  $n$ -exact but not  $(n-1)$ -exact, contradicting Lemma 2.29.

Now suppose  $2^{\text{expo}(x)+1} \neq \text{trunc}(x + 2^{\text{expo}(x)-n}, n)$ . Since  $2^{\text{expo}(x)+1}$  is  $n$ -exact,  $2^{\text{expo}(x)+1} < \text{trunc}(x + 2^{\text{expo}(x)-n}, n)$  by Lemma 2.22. But then by Lemma 2.13,

$$\text{trunc}(x + 2^{\text{expo}(x)-n}, n) \geq 2^{\text{expo}(x)+1} + 2^{\text{expo}(x)+2-n} > x + 2^{\text{expo}(x)-n}.$$

□

**Lemma 2.31.** *If  $n \in \mathbb{Z}$ ,  $n > 1$ , and  $x \in \mathbb{Q}^+$ , then*

$$\text{near}(x, n) = \begin{cases} \text{trunc}(x + 2^{\text{expo}(x)-n}, n - 1) & \text{if } x \text{ is } (n + 1)\text{-exact but not } n\text{-exact} \\ \text{trunc}(x + 2^{\text{expo}(x)-n}, n) & \text{otherwise.} \end{cases}$$

*Proof.* If  $x + 2^{\text{expo}(x)-n} \geq 2^{\text{expo}(x)+1}$ , then by Lemmas 2.19 and 2.30,

$$\text{near}(x, n) = 2^{\text{expo}(x)+1} = \text{trunc}(x + 2^{\text{expo}(x)-n}, n) = \text{trunc}(x + 2^{\text{expo}(x)-n}, n - 1).$$

Thus, we may assume  $x + 2^{\text{expo}(x)-n} < 2^{\text{expo}(x)+1}$ , and it follows from Lemmas 2.21 and 2.26 that

$$\text{expo}(\text{near}(x, n)) = \text{expo}(x + 2^{\text{expo}(x)-n}) = \text{expo}(x).$$

*Case 1.  $x$  is  $n$ -exact*

By Lemma 2.22,  $\text{trunc}(x + 2^{\text{expo}(x)-n}, n) \geq x$ . But since

$$\text{trunc}(x + 2^{\text{expo}(x)-n}, n) \leq x + 2^{\text{expo}(x)-n} < x + 2^{\text{expo}(x)+1-n},$$

Lemma 2.13 yields  $\text{trunc}(x + 2^{\text{expo}(x)-n}, n) \leq x$ ; hence

$$\text{trunc}(x + 2^{\text{expo}(x)-n}, n) = x = \text{near}(x, n).$$

*Case 2.  $x$  is not  $(n + 1)$ -exact*

We have  $\text{near}(x, n) > x - 2^{\text{expo}(x)-n}$ , for otherwise we would have  $\text{near}(x, n) = x - 2^{\text{expo}(x)-n}$  by Lemma 2.26, and since  $\text{near}(x, n)$  is  $(n + 1)$ -exact, so would be

$$\text{near}(x, n) + 2^{\text{expo}(\text{near}(x, n))-n} = x - 2^{\text{expo}(x)-n} + 2^{\text{expo}(\text{near}(x, n))-n} = x.$$

Since  $\text{near}(x, n) \leq x + 2^{\text{expo}(x)-n}$ ,  $\text{near}(x, n) \leq \text{trunc}(x + 2^{\text{expo}(x)-n}, n)$  by Lemma 2.22. But since

$$\text{trunc}(x + 2^{\text{expo}(x)-n}, n) \leq x + 2^{\text{expo}(x)-n} < \text{near}(x, n) + 2^{\text{expo}(x)+1-n},$$

$\text{trunc}(x + 2^{\text{expo}(x)-n}, n) \leq \text{near}(x, n)$ .

*Case 3.  $x$  is  $(n + 1)$ -exact but not  $n$ -exact*

First suppose  $\text{near}(x, n) > x$ . Since  $\text{near}(x, n)$  is  $(n + 1)$ -exact,  $\text{near}(x, n) \geq x + 2^{\text{expo}(x)-n}$ ; hence  $\text{near}(x, n) = x + 2^{\text{expo}(x)-n}$ , and by Lemma 2.29,

$$\text{trunc}(x + 2^{\text{expo}(x)-n}, n - 1) = \text{trunc}(\text{near}(x, n), n - 1) = \text{near}(x, n).$$

Now suppose  $\text{near}(x, n) < x$ . Then  $\text{near}(x, n) < x + 2^{\text{expo}(x)-n}$  implies  $\text{near}(x, n) \leq \text{trunc}(x + 2^{\text{expo}(x)-n}, n - 1)$ . But since

$$\begin{aligned} \text{trunc}(x + 2^{\text{expo}(x)-n}, n - 1) &\leq x + 2^{\text{expo}(x)-n} = x - 2^{\text{expo}(x)-n} + 2^{\text{expo}(x)+1-n} \\ &< \text{near}(x, n) + 2^{\text{expo}(x)+2-n}, \end{aligned}$$

we have  $\text{trunc}(x + 2^{\text{expo}(x)-n}, n - 1) \leq \text{near}(x, n)$ . □

### 3. Multiplication

#### 3.1. The program FPU-MUL

The multiplication algorithm is represented by the program *FPU-MUL*, as listed in Figures 1 and 2. The program is coded in a simple language, consisting of assignments and conditional branches. The primitive operations are logical operations on bit vectors and integer addition and multiplication, the implementation of which is not addressed here.

The algorithm is intended to be implemented with three distinct (integer) multipliers, which operate on the same two  $M$ -bit factors, yielding identical products of either  $2M$  or  $2M - 1$  bits. The output of the first multiplier is manipulated under the assumption that *overflow* occurs, i.e., the product has  $2M$  bits. In parallel, the output of the second multiplier is similarly manipulated under the opposite assumption. Meanwhile, the most significant bit produced by the third multiplier is examined to determine which of the first two results will actually be used, while the other is discarded.

The inputs to this program include two encodings,  $x$  and  $y$ , of the numbers to be multiplied, as well as two specifiers,  $rc$  and  $pc$ , which control the rounding of the product. Irrespective of this rounding, the result is returned in the  $(M, 18)$  format. Thus, the output  $z$  is expected to satisfy

$$\hat{z} = \text{rnd}(\hat{x}\hat{y}, rc, pc).$$

As a notational convenience, the following function gives the position of the *least significant bit* of a  $2M$ -bit integer that has been rounded to a given degree of precision.

**Definition 3.1.** For any precision control specifier  $\pi$ ,  $\text{lsb}(\pi) = 2M - \text{mbits}(\pi)$ .

In addition to computing products, the multiplication hardware performs several auxiliary functions in support of the divide and square root operations. These are specified by the input  $op$ , the value of which may be any of the symbols *OP-MUL*, *OP-DIV*, *OP-SQRT*, *OP-LAST*, and *OP-BACK*.

Basic floating point multiplication is performed in the case  $op = \text{OP-MUL}$ : the inputs  $x$  and  $y$  are simply multiplied and rounded according to the specifiers  $pc$  and  $rc$ , and the IEEE compliant result is returned as the output  $z$ , as described by Theorem 1. The same holds for  $op = \text{OP-DIV}$  and  $op = \text{OP-SQRT}$ , but an additional output  $r$  is returned in these cases: for *OP-DIV*,  $\hat{r}$  is an approximation of  $2 - \hat{x}\hat{y}$ ; for *OP-SQRT*,  $\hat{r}$  is an approximation of  $(3 - \hat{x}\hat{y})/2$ . The errors of these approximations are given by Lemma 3.5.

When *FPU-MUL* is called by division or square root,  $pc$  is always *PC-\**, indicating the internal format  $(M, 18)$ . However, on the final iteration of either of these operations, signaled by *OP-LAST*, the product is rounded to a lower precision, as determined by the input  $lastpc$ . This behavior is described formally by Lemma 3.7.

Finally, the symbol *OP-BACK* indicates a *back multiplication* to determine whether the product previously computed by *OP-LAST* is an overestimate or an underestimate of the exact value sought. The value given by the input  $d$  is subtracted from the product of  $x$  and  $y$ . In the case of division,  $x$  is the denominator,  $y$  is the approximate quotient, and  $d$  is the numerator; in the square root case, both  $x$  and  $y$  are the approximate square root and  $d$  is the radicand. In both cases, the results of the comparison are given by the outputs  $z$  and  $inexact$ , as stated in Lemma 3.8.

Thus, our analysis will be based on an execution of

$$FPU-MUL(op, pc, lastpc, rc, x, y, z, r, d, inexact),$$



**Program**  $FPU-MUL(op, pc, lastpc, rc, x, y, z, r, d, inexact)$ :

```

sign ← get-sign(x) ^ get-sign(y);
man-unrounded ← get-man(x) · get-man(y);
overflow ← man-unrounded[2M - 1];
if man-unrounded[lsb(pc) - 3 : 0] = 0
    then sticky-no-overflow ← 0
    else sticky-no-overflow ← 1;
sticky-with-overflow ← sticky-no-overflow | man-unrounded[lsb(pc) - 2];
inexact-no-overflow ← sticky-with-overflow;
inexact-with-overflow ← inexact-no-overflow | man-unrounded[lsb(pc) - 1];
if op = OP-BACK
    then if overflow = 1
        then inexact ← inexact-with-overflow
        else inexact ← inexact-no-overflow;
if op = OP-BACK then
    rconst-with-overflow ← compl(2M get-man(d), 2M)
else if op = OP-LAST then
    rconst-with-overflow ← 2lsb(lastpc)-2
else if rc = RC-NEAR then
    rconst-with-overflow ← 2lsb(pc)-1
else if (sign = 1 ∧ rc = RC-NEG) ∨ (sign = 0 ∧ rc = RC-POS) then
    rconst-with-overflow ← 2lsb(pc) - 1
else rconst-with-overflow ← 0;
rconst-no-overflow ← shr(rconst-with-overflow, 0, 2M);
if op = OP-BACK
    then {add-with-overflow ← (man-unrounded + rconst-with-overflow + 1)[2M : 0];
        add-no-overflow ← (man-unrounded + rconst-no-overflow + 1)[2M - 1 : 0]}
    else {add-with-overflow ← (man-unrounded + rconst-with-overflow)[2M : 0];
        add-no-overflow ← (man-unrounded + rconst-no-overflow)[2M - 1 : 0]};
round-carryout-no-overflow ← add-no-overflow[2M - 1];
round-carryout-with-overflow ← add-with-overflow[2M];
if op = OP-LAST
    then {trunc-with-overflow ← 22M - 2lsb(lastpc)-1;
        trunc-no-overflow ← 22M - 2lsb(lastpc)-2}
    else {trunc-with-overflow ← 22M - 2lsb(pc);
        trunc-no-overflow ← 22M - 2lsb(pc)-1};

```

Figure 1:  $FPU-MUL$

```

if rc = RC-NEAR ∧ sticky-no-overflow = 0 ∧ add-no-overflow[lsb(pc) - 2] = 0
  then man-rounded-no-overflow
    ← (22M-2round-carryout-no-overflow | add-no-overflow)
      & ((22M - 1 - 2lsb(pc)-1) & trunc-no-overflow)
  else man-rounded-no-overflow
    ← (22M-2round-carryout-no-overflow | add-no-overflow)
      & trunc-no-overflow;
if rc = RC-NEAR ∧ sticky-with-overflow = 0 ∧ add-with-overflow[lsb(pc) - 1] = 0
  then man-rounded-with-overflow
    ← (22M-1round-carryout-with-overflow | add-with-overflow)
      & ((22M - 1 - 2lsb(pc)) & trunc-with-overflow);
  else man-rounded-with-overflow
    ← (22M-1round-carryout-with-overflow | add-with-overflow)
      & trunc-with-overflow;
exp-unrounded ← (get-expo(x) + get-expo(y) + 217 + 1)[17 : 0];
exp-rounded-with-overflow
  ← (exp-unrounded + round-carryout-with-overflow + 1)[17 : 0];
exp-rounded-no-overflow ← (exp-unrounded + round-carryout-no-overflow)[17 : 0];
if get-man(x) = 0 then
  z ← (sign, 0, get-expo(x))
else if get-man(y) = 0 then
  z ← (sign, 0, get-expo(y))
else if overflow = 1 then
  z ← (sign, man-rounded-with-overflow[2M - 1 : M], exp-rounded-with-overflow)
  else z ← (sign, man-rounded-no-overflow[2M - 2 : M - 1], exp-rounded-no-overflow);
if op = OP-DIV then
  if overflow = 1 then
    r ← (0, compl(man-unrounded, 2M)[2M - 2 : M - 1], 217 - 2)
  else if round-carryout-no-overflow = 0 then
    r ← (0, compl(man-unrounded, 2M)[2M - 1 : M], 217 - 1)
  else r ← (0, 2M - 1, 217 - 2)
else if op = OP-SQRT then
  if overflow = 1 then
    r ← (0, (compl(man-unrounded, 2M) | 22M-1)[2M - 1 : M], 217 - 2)
  else if round-carryout-no-overflow = 0 then
    r ← (0, shr(compl(man-unrounded, 2M)[2M - 2 : 0], 1, 2M)[2M - 1 : M], 217 - 1)
  else r ← (0, 2M - 1, 217 - 2)

```

Figure 2: FPU-MUL (continued)

under the following assumptions regarding the inputs.

- (a)  $op \in \{OP-MUL, OP-DIV, OP-SQRT, OP-LAST, OP-BACK\}$ ;
- (b)  $pc$  is a precision control specifier;
- (c) if  $op = OP-LAST$ , then  $lastpc$  is an external precision control specifier;
- (d)  $rc$  is a rounding control specifier;
- (e)  $x$  and  $y$  are normal encodings;
- (f) if  $op = OP-BACK$ , then  $d$  is a normal encoding.

### 3.2. Basic results

For convenience, we introduce several auxiliary variables. First, we define

$$sticky = \begin{cases} sticky-with-overflow & \text{if overflow} = 1 \\ sticky-no-overflow & \text{if overflow} = 0. \end{cases}$$

Each of the variables  $rconst$ ,  $add$ ,  $round-carryout$ ,  $trunc$ ,  $man-rounded$ , and  $expo-rounded$  is defined in the analogous manner. We also define

$$P = \begin{cases} 2M & \text{if overflow} = 1 \\ 2M - 1 & \text{if overflow} = 0, \end{cases}$$

$$\mu = \text{mbits}(pc),$$

and

$$trunc' = \begin{cases} trunc, & \text{if } rc \neq RC-NEAR \text{ or } sticky = 1 \text{ or } add[P - \mu - 1] = 1 \\ trunc \ \& \ (2^{2M} - 1 - 2^{P-\mu}), & \text{otherwise.} \end{cases}$$

#### Lemma 3.1.

- (a)  $\text{sig}(\text{man-unrounded}) = \text{sig}(\hat{x})\text{sig}(\hat{y})/2^{\text{overflow}}$ ;
- (b)  $\text{expo}(\text{man-unrounded}) = P - 1$ ;
- (c)  $\text{sig}(\hat{x}\hat{y}) = \text{sig}(\text{man-unrounded})$ ;
- (d)  $\text{expo}(\hat{x}\hat{y}) = \text{expo}(\hat{x}) + \text{expo}(\hat{y}) + \text{overflow}$ .

*Proof.* Since  $x$  and  $y$  are normal encodings,

$$2^{2M-2} \leq \text{man-unrounded} = \text{get-man}(x) \cdot \text{get-man}(y) < 2^{2M},$$

and (b) follows from Lemma 2.2.

By Lemma 2.15,

$$\begin{aligned} \text{man-unrounded} &= 2^{M-1} \text{sig}(\hat{x}) 2^{M-1} \text{sig}(\hat{y}) \\ &= \text{sig}(\hat{x}) \text{sig}(\hat{y}) 2^{-\text{overflow}} 2^{2M-2+\text{overflow}} \\ &= \text{sig}(\hat{x}) \text{sig}(\hat{y}) 2^{-\text{overflow}} 2^{\text{expo}(\text{man-unrounded})}, \end{aligned}$$

which implies (a).

To derive (c) and (d), we need only observe that

$$\begin{aligned} \hat{x}\hat{y} &= \text{sgn}(\hat{x}) \text{sig}(\hat{x}) 2^{\text{expo}(\hat{x})} \text{sgn}(\hat{y}) \text{sig}(\hat{y}) 2^{\text{expo}(\hat{y})} \\ &= \text{sgn}(\hat{x}\hat{y}) \left[ \text{sig}(\hat{x}) \text{sig}(\hat{y}) / 2^{\text{overflow}} \right] 2^{\text{expo}(\hat{x}) + \text{expo}(\hat{y}) + \text{overflow}}. \end{aligned}$$

**Lemma 3.2.**

- (a) sticky = 0 iff man-unrounded is  $(\mu + 1)$ -exact;
- (b) inexact = 0 iff man-unrounded is  $\mu$ -exact.

*Proof.* We have sticky-no-overflow = 0  $\Leftrightarrow 2^{\text{lsb}(\text{pc})-2}$  divides man-unrounded, and

sticky-with-overflow = 0

- $\Leftrightarrow 2^{\text{lsb}(\text{pc})-2}$  divides man-unrounded and man-unrounded[lsb(pc) - 2] = 0
- $\Leftrightarrow 2^{\text{lsb}(\text{pc})-2}$  divides man-unrounded and 2 divides man-unrounded/ $2^{\text{lsb}(\text{pc})-2}$
- $\Leftrightarrow 2^{\text{lsb}(\text{pc})-1}$  divides man-unrounded.

Thus, sticky = 0 iff  $2^{P-(\mu+1)}$  divides man-unrounded, and (a) follows from Lemma 2.14.

Similarly, it may be shown that inexact = 0 iff  $2^{P-\mu}$  divides man-unrounded, which implies (b). □

**Lemma 3.3.**

- (a) man-rounded =  $(2^{P-1}\text{round-carryout} \mid \text{add} \ \& \ \text{trunc}')$ ;
- (b) man-rounded[ $P - 1$ ] = 1;
- (c) expo(man-rounded)  $\leq$  expo(add) =  $P - 1 + \text{round-carryout}$ ;
- (d) man-rounded is divisible by  $2^{P-M}$ .

*Proof.* (a) In all cases,

$$\text{man-rounded} = (2^{P-1}\text{round-carryout} \mid \text{add}) \ \& \ \text{trunc}'$$

and  $\text{trunc}'[P - 1] = 1$ . Thus, by Lemmas 2.7 and 2.11,

$$\begin{aligned} \text{man-rounded} &= (2^{P-1}\text{round-carryout} \ \& \ \text{trunc}') \mid (\text{add} \ \& \ \text{trunc}') \\ &= 2^{P-1}\text{round-carryout} \mid (\text{add} \ \& \ \text{trunc}') \end{aligned}$$

- (b) By Lemma 2.10, we may assume round-carryout = 0 and hence

$$\text{man-rounded}[P - 1] = \text{add}[P - 1].$$

Note that

$$\text{add} = \begin{cases} \text{rem}(\text{man-unrounded} + \text{rconst} + 1, 2^{P+1}) & \text{if op} = \text{OP-BACK} \\ \text{rem}(\text{man-unrounded} + \text{rconst}, 2^{P+1}) & \text{otherwise,} \end{cases}$$

and that since

$$\text{man-unrounded} + \text{rconst} + 1 \leq (2^P - 1) + (2^P - 1) + 1 < 2^{P+1},$$

we have

$$2^{P-1} \leq \text{man-unrounded} \leq \text{add} < 2^{P+1}.$$

But since round-carryout = add[ $P$ ] = 0, Lemma 2.2 implies add <  $2^P$  and hence add[ $P - 1$ ] = 1.

- (c) If round-carryout = 0, then

$$\text{man-rounded} = \text{add} \ \& \ \text{trunc}' \leq \text{add} < 2^P,$$

by Lemma 2.9, and  $\text{man-rounded}[P - 1] = 1$  implies  $\text{man-rounded} \geq 2^{P-1}$ ; hence

$$\text{expo}(\text{man-rounded}) = \text{expo}(\text{add}) = P - 1.$$

On the other hand, if  $\text{round-carryout} = \text{add}[P] = 1$ , then  $\text{expo}(\text{add}) = P$ , while  $\text{man-rounded} < 2^{P+1}$  by Lemma 2.8; hence  $\text{expo}(\text{man-rounded}) \leq P$ .

(d) Since  $2^{P-M}$  divides  $\text{trunc}$ , the result follows from Lemmas 2.9 and 2.8. □

**Lemma 3.4.**  *$z$  is a normal encoding and*

- (a)  $\text{sgn}(\hat{z}) = \text{sgn}(\hat{x}\hat{y})$ ;
- (b)  $\text{sig}(\hat{z}) = \text{rem}(\text{man-rounded}, 2^P)/2^{P-1}$ ;
- (c)  $\text{rem}(\text{expo}(\hat{z}), 2^{18}) = \text{rem}(\text{expo}(\hat{x}\hat{y}) + \text{round-carryout}, 2^{18})$ .

*Proof.* First, observe that

$$z = (\text{sign}, \text{man-rounded}[P - 1 : P - M], \text{exp-rounded}).$$

Let  $\rho = \text{rem}(\text{man-rounded}, 2^P)$ . By Lemma 2.3,

$$\rho[P - 1] = \text{man-rounded}[P - 1] = 1;$$

hence  $\text{expo}(\rho) = P - 1$ . Since  $\text{man-rounded}$  is divisible by  $2^{P-M}$ , so is  $\rho$ . Thus, by Lemmas 2.4 and 2.14,

$$\text{get-man}(z) = \text{man-rounded}[P - 1 : P - M] = \rho[P - 1 : P - M] = \rho/2^{P-M}.$$

It follows that

$$\text{expo}(\text{get-man}(z)) = \text{expo}(\rho) - (P - M) = (P - 1) - (P - M) = M - 1.$$

Since

$$\text{get-expo}(z) = \text{exp-rounded} = \text{rem}(\text{exp-unrounded} + \text{round-carryout} + \text{overflow}, 2^{18}),$$

$0 < \text{get-expo}(z) < 2^{18}$ , and hence  $z$  is a normal encoding. The proof is completed by applying Lemma 2.15.

(a)  $\text{sgn}(\hat{z}) = (-1)^{\text{sign}}$ ; hence  $\text{sgn}(\hat{z}) = 1 \Leftrightarrow \text{sign} = 0 \Leftrightarrow \text{get-sign}(x) = \text{get-sign}(y) \Leftrightarrow \text{sgn}(\hat{x}) = \text{sgn}(\hat{y}) \Leftrightarrow \text{sgn}(\hat{x}\hat{y}) = 1$ .

(b)  $\text{sig}(\hat{z}) = \text{get-man}(z)/2^{M-1} = (\rho/2^{P-M})/2^{M-1} = \rho/2^{P-1}$ .

(c)  $\text{expo}(\hat{z}) = \text{get-expo}(z) - (2^{17} - 1)$ , where

$$\begin{aligned} &\text{get-expo}(z) \\ &= \text{rem}(\text{exp-unrounded} + \text{round-carryout} + \text{overflow}, 2^{18}) \\ &= \text{rem}(\text{get-expo}(x) + \text{get-expo}(y) + 2^{17} + 1 + \text{round-carryout} + \text{overflow}, 2^{18}) \\ &= \text{rem}(\text{expo}(\hat{x}) + \text{expo}(\hat{y}) + 2^{18} - 2 + 2^{17} + 1 + \text{round-carryout} + \text{overflow}, 2^{18}) \\ &= \text{rem}(\text{expo}(\hat{x}) + \text{expo}(\hat{y}) + \text{overflow} + 2^{17} - 1 + \text{round-carryout}, 2^{18}) \\ &= \text{rem}(\text{expo}(\hat{x}\hat{y}) + 2^{17} - 1 + \text{round-carryout}, 2^{18}). \end{aligned} \quad \square$$

### 3.3. The operations OP-MUL, OP-DIV, and OP-SQRT

This is our statement of IEEE compliance for multiplication.

**Theorem 1.** *Assume that  $\text{op} \in \{\text{OP-MUL}, \text{OP-DIV}, \text{OP-SQRT}\}$ ,  $\text{rc}$  is a rounding control specifier,  $\text{pc}$  is a precision control specifier, and  $x$  and  $y$  are normal encodings. If  $\text{rnd}(\hat{x}\hat{y}, \text{rc}, \text{pc})$  is representable, then  $\hat{z} = \text{rnd}(\hat{x}\hat{y}, \text{rc}, \text{pc})$ .*

*Proof.* Let

$$rc' = \begin{cases} \text{RC-NEG} & \text{if } rc = \text{RC-POS} \\ \text{RC-POS} & \text{if } rc = \text{RC-NEG} \\ rc & \text{otherwise.} \end{cases}$$

Then  $\text{rnd}(-\hat{x}\hat{y}, rc, pc) = -\text{rnd}(\hat{x}\hat{y}, rc', pc)$ . Also, by inspection of the code that defines *FPU-MUL*, it is easy to see that replacing either “get – sign(*x*)” or “get – sign(*y*)” by its complement and *rc* by *rc'* has the effect of negating  $\hat{z}$ . It follows that it suffices to prove the theorem under the assumptions  $\hat{x} > 0$  and  $\hat{y} > 0$ , which imply that  $\text{sign} = 0$ .

Note that (under these assumptions)

$$\text{rconst} = \begin{cases} 2^{P-\mu-1} & \text{if } rc = \text{RC-NEAR} \\ 2^{P-\mu} - 1 & \text{if } rc = \text{RC-POS} \\ 0 & \text{otherwise.} \end{cases}$$

In all cases,  $\text{rconst} < 2^P$ . Since  $\text{man-unrounded} < 2^P$  as well,

$$\text{add} = \text{rem}(\text{man-unrounded} + \text{rconst}, 2^{P+1}) = \text{man-unrounded} + \text{rconst}.$$

If  $rc = \text{RC-NEAR}$  and  $\text{sticky} = \text{add}[P - \mu - 1] = 0$ , then by Lemma 2.12,

$$\text{trunc}' = (2^{2M} - 2^{P-\mu}) \& (2^{2M} - 1 - 2^{P-\mu}) = (2^{2M} - 2^{P-\mu+1}),$$

and otherwise

$$\text{trunc}' = (2^{2M} - 2^{P-\mu}).$$

We shall show that

$$\text{rem}(\text{man-rounded}, 2^P) = \text{rnd}(\text{man-unrounded}, rc, pc)2^{-\text{round-carryout}},$$

by considering the following cases.

*Case 1.*  $\text{round-carryout} = 0$

Since  $\text{man-rounded} < 2^P$  by Lemma 3.3, we must show

$$\text{man-rounded} = \text{rnd}(\text{man-unrounded}, rc, pc).$$

*Subcase 1.1.*  $rc = \text{RC-NEAR}$

First suppose  $\text{sticky} = \text{add}[P - \mu - 1] = 0$ . Then Lemma 2.3 implies

$$\text{man-unrounded}[P - \mu - 1] = 1,$$

and by Lemmas 3.2 and 2.14,  $\text{man-unrounded}$  is  $(\mu + 1)$ -exact but not  $\mu$ -exact. Thus, by Lemmas 3.3, 2.24, and 2.31,

$$\begin{aligned} \text{man-rounded} &= (\text{man-unrounded} + 2^{P-\mu-1}) \& (2^{2M} - 2^{P-\mu+1}) \\ &= \text{trunc}(\text{man-unrounded} + 2^{P-\mu-1}, \mu - 1) \\ &= \text{near}(\text{man-unrounded}, \mu) \\ &= \text{rnd}(\text{man-unrounded}, rc, pc). \end{aligned}$$

In the remaining case,  $\text{man-unrounded}$  is either  $\mu$ -exact or not  $(\mu + 1)$ -exact, and the same

three lemmas yield

$$\begin{aligned} \text{man-rounded} &= (\text{man-unrounded} + 2^{P-\mu-1}) \& (2^{2M} - 2^{P-\mu}) \\ &= \text{trunc}(\text{man-unrounded} + 2^{P-\mu-1}, \mu) \\ &= \text{near}(\text{man-unrounded}, \mu) \\ &= \text{rnd}(\text{man-unrounded}, \text{rc}, \text{pc}). \end{aligned}$$

Subcase 1.2. rc = RC-POS

By Lemmas 2.24 and 2.25,

$$\begin{aligned} \text{man-rounded} &= (\text{man-unrounded} + 2^{P-\mu} - 1) \& (2^{2M} - 2^{P-\mu}) \\ &= \text{trunc}(\text{man-unrounded} + 2^{P-\mu} - 1, \mu) \\ &= \text{away}(\text{man-unrounded}, \mu) \\ &= \text{rnd}(\text{man-unrounded}, \text{rc}, \text{pc}). \end{aligned}$$

Subcase 1.3. rc = RC-CHOP or rc = RC-NEG

By Lemma 2.24,

$$\begin{aligned} \text{man-rounded} &= \text{man-unrounded} \& (2^{2M} - 2^{P-\mu}) \\ &= \text{trunc}(\text{man-unrounded}, \mu) \\ &= \text{rnd}(\text{man-unrounded}, \text{rc}, \text{pc}). \end{aligned}$$

Case 2. round-carryout = 1

In this case,

$$2^P \leq \text{add} = \text{man-unrounded} + \text{rconst} < 2^P + \text{rconst},$$

which implies

$$0 \leq \text{rem}(\text{add}, 2^P) < \text{rconst} < 2^{P-\mu}.$$

By Lemmas 3.3, 2.9, and 2.8,

$$\begin{aligned} \text{rem}(\text{man-rounded}, 2^P) &= \text{rem}(2^{P-1} \mid (\text{add} \& \text{trunc}'), 2^P) \\ &= 2^{P-1} \mid (\text{rem}(\text{add}, 2^P) \& \text{trunc}') \\ &= 2^{P-1} \mid (\text{rem}(\text{add}, 2^P) \& \text{rem}(\text{trunc}', 2^{P-\mu})) \\ &= 2^{P-1} \mid (\text{rem}(\text{add}, 2^P) \& 0) \\ &= 2^{P-1}. \end{aligned}$$

Thus, it suffices to show that  $\text{rnd}(\text{man-unrounded}, \text{rc}, \text{pc}) = 2^P$ .

Subcase 2.1. rc = RC-NEAR

Since

$$\text{man-unrounded} + 2^{P-1-\mu} = \text{man-unrounded} + \text{rconst} \geq 2^P,$$

$\text{near}(\text{man-unrounded}, \mu) = 2^P$  by Lemma 2.30.

Subcase 2.2. rc = RC-POS

Let  $a = 2^P - 2^{P-\mu}$ . Then

$$\text{man-unrounded} \geq 2^P - \text{rconst} = 2^P - 2^{P-\mu} + 1 > a,$$

and since  $a$  is  $\mu$ -exact,

$$\text{away}(\text{man-unrounded}, \mu) \geq a + 2^{\text{expo}(a)+1-\mu} = a + 2^{P-\mu} = 2^P,$$

which implies  $\text{away}(\text{man-unrounded}, \mu) = 2^P$ .

*Subcase 2.3.*  $\text{rc} = \text{RC-CHOP}$  or  $\text{rc} = \text{RC-NEG}$

This case is precluded by our earlier observation that  $0 < \text{const}$ .

The proof is completed by applying Lemmas 3.4 and 3.1, which yield

$$\text{sgn}(\hat{z}) = \text{sgn}(\hat{x}\hat{y}) = 1,$$

$$\begin{aligned} \text{sig}(\hat{z}) &= \text{rnd}(\text{man-unrounded}, \text{rc}, \text{pc})2^{-\text{round-carryout}-P+1} \\ &= \text{rnd}(\text{sig}(\hat{x}\hat{y}), \text{rc}, \text{pc})2^{-\text{round-carryout}}, \end{aligned}$$

and for some  $k \in \mathbb{Z}$ ,

$$\text{expo}(\hat{z}) = \text{expo}(\hat{x}\hat{y}) + \text{round-carryout} + 2^{18}k.$$

Thus,

$$\hat{z} = \text{rnd}(\text{sig}(\hat{x}\hat{y}), \text{rc}, \text{pc})2^{\text{expo}(\hat{x}\hat{y})+2^{18}k} = \text{rnd}(\hat{x}\hat{y}, \text{rc}, \text{pc})2^{2^{18}k}.$$

But since  $\text{rnd}(\hat{x}\hat{y}, \text{rc}, \text{pc})$  is representable, i.e.,  $1 - 2^{-17} \leq \text{expo}(\text{rnd}(\hat{x}\hat{y}, \text{rc}, \text{pc})) \leq 2^{17}$ , and the same is true of  $\hat{z}$ ,

$$|2^{18}k| = |\text{expo}(\hat{z}) - \text{expo}(\text{rnd}(\hat{x}\hat{y}, \text{rc}, \text{pc}))| < 2^{18},$$

and hence  $k = 0$ . □

In the  $\text{OP-DIV}$  and  $\text{OP-SQRT}$  cases, an additional value is returned.

**Lemma 3.5.** *Let  $\text{op} \in \{\text{OP-DIV}, \text{OP-SQRT}\}$ ,  $\text{pc} = \text{PC-}^*$ , and  $\text{rc} = \text{RC-NEAR}$ . Assume that  $x$  and  $y$  are normal encodings,  $3/2 < \text{sig}(\hat{x})\text{sig}(\hat{y}) < 3$ , and  $|1 - \hat{x}\hat{y}| < 1/8$ . Then*

- (a)  $r$  is a normal encoding;
- (b)  $\hat{r} < 1 \Leftrightarrow \hat{z} \geq 1$ ;
- (c) if  $\text{op} = \text{OP-DIV}$ , then  $2 - \hat{x}\hat{y} - 2^{1-M} \leq \hat{r} < 2 - \hat{x}\hat{y}$ ;
- (d) if  $\text{op} = \text{OP-SQRT}$ , then  $(3 - \hat{x}\hat{y})/2 - 2^{1-M} \leq \hat{r} < (3 - \hat{x}\hat{y})/2$ .

*Proof.* First note that the hypothesis implies that  $\text{expo}(\hat{x}\hat{y})$  is either 0 or  $-1$ , and it follows from Lemma 3.4 that

$$\text{expo}(\hat{z}) = \text{expo}(\hat{x}\hat{y}) + \text{round-carryout}.$$

We consider the following cases.

*Case 1.*  $\text{overflow} = 1$

In this case,  $\text{expo}(\text{man-unrounded}) = 2M - 1$ , but by Lemma 3.1,

$$\text{man-unrounded} = \text{sig}(\hat{x}\hat{y})2^{2M-1} = \text{sig}(\hat{x})\text{sig}(\hat{y})2^{2M-2} < 3 \cdot 2^{2M-2}$$

and hence

$$\text{add} = \text{man-unrounded} + 2^{M-1} < 3 \cdot 2^{2M-2} + 2^{M-1} < 2^{2M}$$

and  $\text{round-carryout} = 0$ . We have  $\text{expo}(\hat{z}) = \text{expo}(\hat{x}\hat{y}) = 0$ , for otherwise

$$\hat{x}\hat{y} = \text{sig}(\hat{x}\hat{y})/2 = \text{sig}(\hat{x})\text{sig}(\hat{y})/4 < 3/4,$$



contradicting  $|1 - \hat{x}\hat{y}| < 1/8$ . Thus,

$$\hat{x}\hat{y} = \text{sig}(\hat{x}\hat{y}) = \text{sig}(\text{man-unrounded}) = \text{man-unrounded}/2^{2M-1}.$$

Also note that

$$\begin{aligned} \text{comp1}(\text{man-unrounded}, 2M) &= 2^{2M} - \text{man-unrounded} - 1 \\ &\leq 2^{2M} - 2^{2M-1} - 1 < 2^{2M-1}. \end{aligned}$$

*Subcase 1.1.*  $\text{op} = \text{OP-DIV}$

$$\begin{aligned} \text{get-man}(r) &= \text{comp1}(\text{man-unrounded}, 2M)[2M - 2 : M - 1] \\ &= \lfloor (2^{2M} - \text{man-unrounded} - 1)/2^{M-1} \rfloor \\ &= 2^{M+1} + \lfloor -(\text{man-unrounded} + 1)/2^{M-1} \rfloor \\ &= 2^{M+1} - \lfloor \text{man-unrounded}/2^{M-1} \rfloor - 1. \end{aligned}$$

But

$$\lfloor \text{man-unrounded}/2^{M-1} \rfloor \leq \text{man-unrounded}/2^{M-1} = 2^M \hat{x}\hat{y}$$

and

$$\lfloor \text{man-unrounded}/2^{M-1} \rfloor > \text{man-unrounded}/2^{M-1} - 1 = 2^M \hat{x}\hat{y} - 1;$$

hence

$$2^{M-1} \leq 2^{M+1} - 2^M \hat{x}\hat{y} - 1 \leq \text{get-man}(r) < 2^{M+1} - 2^M \hat{x}\hat{y} \leq 2^M$$

and  $r$  is normal. Since  $\text{expo}(\hat{r}) = 2^{17} - 2 - (2^{17} - 1) = -1$ ,  $\hat{r} < 1 \leq \hat{z}$  and

$$2 - \hat{x}\hat{y} - 2^{-M} \leq \hat{r} = 2^{-M} \text{get-man}(r) < 2 - \hat{x}\hat{y}.$$

*Subcase 1.2.*  $\text{op} = \text{OP-SQRT}$

By Lemmas 2.2 and 2.11,

$$\begin{aligned} \text{comp1}(\text{man-unrounded}, 2M) \mid 2^{2M-1} &= \text{comp1}(\text{man-unrounded}, 2M) + 2^{2M-1} \\ &= 2^{2M} + 2^{2M-1} - \text{man-unrounded} - 1 \\ &< 2^{2M}; \end{aligned}$$

hence

$$\begin{aligned} \text{get-man}(r) &= (\text{comp1}(\text{man-unrounded}, 2M) \mid 2^{2M-1})[2M - 1 : M] \\ &= (2^{2M} + 2^{2M-1} - \text{man-unrounded} - 1)[2M - 1 : M] \\ &= \lfloor (2^{2M} + 2^{2M-1} - \text{man-unrounded} - 1)/2^M \rfloor \\ &= 2^M + 2^{M-1} + \lfloor -(\text{man-unrounded} + 1)/2^M \rfloor \\ &= 3 \cdot 2^{M-1} - \lfloor \text{man-unrounded}/2^M \rfloor - 1. \end{aligned}$$

But

$$\lfloor \text{man-unrounded}/2^M \rfloor \leq \text{man-unrounded}/2^M = 2^{M-1} \hat{x}\hat{y}$$

and

$$\lfloor \text{man-unrounded}/2^M \rfloor > \text{man-unrounded}/2^M - 1 = 2^{M-1} \hat{x}\hat{y} - 1,$$

implying

$$2^{M-1} \leq 2^{M-1} (3 - \hat{x}\hat{y}) - 1 \leq \text{get-man}(r) < 2^{M-1} (3 - \hat{x}\hat{y}) \leq 2^M;$$

hence  $r$  is normal. Again,  $\text{expo}(\hat{r}) = -1$  and  $\hat{r} < 1 \leq \hat{z}$ . Thus

$$(3 - \hat{x}\hat{y})/2 - 2^{-M} \leq \hat{r} = \text{get-man}(r)/2^M < (3 - \hat{x}\hat{y})/2.$$

*Case 2. overflow = 0*

In this case,  $\text{expo}(\text{man-unrounded}) = 2M - 2$ , and  $\text{expo}(\hat{x}\hat{y}) = -1$ , for otherwise

$$\hat{x}\hat{y} = \text{sig}(\hat{x}\hat{y}) = \text{sig}(\hat{x})\text{sig}(\hat{y}) > 3/2,$$

contradicting  $|1 - \hat{x}\hat{y}| < 1/8$ . Thus

$$\hat{x}\hat{y} = \text{sig}(\hat{x}\hat{y})/2 = \text{sig}(\text{man-unrounded})/2 = \text{man-unrounded}/2^{2M-1}.$$

*Subcase 2.1. round-carryout = 0*

Since  $\text{expo}(\hat{z}) = \text{expo}(\hat{x}\hat{y}) = -1$ ,  $\hat{z} < 1$ .

*Subcase 2.1.1. op = OP-DIV*

$$\begin{aligned} \text{get-man}(r) &= (2^{2M} - \text{man-unrounded} - 1)[2M - 1 : M] \\ &= \lfloor (2^{2M} - \text{man-unrounded} - 1)/2^M \rfloor \\ &= 2^M + \lfloor -(\text{man-unrounded} + 1)/2^M \rfloor \\ &= 2^M - \lfloor \text{man-unrounded}/2^M \rfloor - 1. \end{aligned}$$

In this case,

$$2^{M-1}\hat{x}\hat{y} - 1 < \lfloor \text{man-unrounded}/2^M \rfloor \leq 2^{M-1}\hat{x}\hat{y}$$

and

$$2^{M-1} - 1 < 2^M - 2^{M-1}\hat{x}\hat{y} - 1 \leq \text{get-man}(r) < 2^M - 2^{M-1}\hat{x}\hat{y} < 2^M.$$

Since  $\text{expo}(\hat{r}) = (2^{17} - 1) - (2^{17} - 1) = 0$ ,  $\hat{r} \geq 1 > \hat{z}$ , and

$$2 - \hat{x}\hat{y} - 2^{1-M} \leq \hat{r} = \text{get-man}(r)/2^{M-1} < 2 - \hat{x}\hat{y}.$$

*Subcase 2.1.2. op = OP-SQRT*

Note that

$$\text{compl}(\text{man-unrounded}, 2M) = 2^{2M} - \text{man-unrounded} - 1 \geq 2^{2M} - 2^{2M-1} = 2^{2M-1},$$

while  $\text{compl}(\text{man-unrounded}, 2M) < 2^{2M}$ ; hence

$$\begin{aligned} \text{rem}(\text{compl}(\text{man-unrounded}, 2^{2M}), 2^{2M-1}) \\ &= \text{compl}(\text{man-unrounded}, 2M) - 2^{2M-1} \\ &= 2^{2M-1} - \text{man-unrounded} - 1. \end{aligned}$$

Therefore, applying Lemma 2.11, we have

$$\begin{aligned}
 \text{get-man}(r) &= \text{shr}(\text{compl}(\text{man-unrounded}, 2M)[2M - 2 : 0], 1, 2M)[2M - 1 : M] \\
 &= \text{shr}(\text{rem}(\text{compl}(\text{man-unrounded}, 2M), 2^{2M-1}), 1, 2M)[2M - 1 : M] \\
 &= \text{shr}(2^{2M-1} - \text{man-unrounded} - 1, 1, 2M)[2M - 1 : M] \\
 &= (2^{2M-1} + \lfloor (2^{2M-1} - \text{man-unrounded} - 1)/2 \rfloor)[2M - 1 : M] \\
 &= \lfloor (2^{2M-1} + \lfloor (2^{2M-1} - \text{man-unrounded} - 1)/2 \rfloor) / 2^M \rfloor \\
 &= 2^{M-1} + \lfloor \lfloor (2^{2M-1} - \text{man-unrounded} - 1)/2 \rfloor / 2^M \rfloor \\
 &= 2^{M-1} + \lfloor (2^{2M-1} - \text{man-unrounded} - 1) / 2^{M+1} \rfloor \\
 &= 2^{M-1} + 2^{M-2} + \lfloor -(\text{man-unrounded} + 1) / 2^{M+1} \rfloor \\
 &= 3 \cdot 2^{M-2} - \lfloor \text{man-unrounded} / 2^{M+1} \rfloor - 1.
 \end{aligned}$$

But

$$2^{M-2} \hat{x} \hat{y} - 1 < \lfloor \text{man-unrounded} / 2^{M+1} \rfloor \leq 2^{M-2} \hat{x} \hat{y};$$

hence

$$2^{M-1} - 1 < 2^{M-2}(3 - \hat{x} \hat{y}) - 1 \leq \text{get-man}(r) < 2^{M-2}(3 - \hat{x} \hat{y}) < 2^M.$$

Again,  $\text{expo}(\hat{r}) = 0$ ,  $\hat{r} \geq 1 > \hat{z}$ , and

$$(3 - \hat{x} \hat{y})/2 - 2^{1-M} \leq \hat{r} = \text{get-man}(r)/2^{M-1} < (3 - \hat{x} \hat{y})/2.$$

*Subcase 2.2.*  $\text{round-carryout} = 1$

In this case,  $\text{get-man}(r) = 2^M - 1$  and  $\hat{r} = 1 - 2^{-M} < 1$ , while  $\text{expo}(\hat{z}) = \text{expo}(\hat{x} \hat{y}) + 1 = 0$ , so  $\hat{z} \geq 1$ . Since  $\text{add} = \text{man-unrounded} + 2^{M-2} \geq 2^{2M-1}$ , we have

$$2^{2M-1} - 2^{M-2} \leq \text{man-unrounded} < 2^{2M-1}$$

and hence

$$1 - 2^{-1-M} \leq \hat{x} \hat{y} < 1,$$

which implies

$$2 - \hat{x} \hat{y} - (2^{-M} + 2^{-1-M}) \leq \hat{r} < 2 - \hat{x} \hat{y} - 2^{-M}$$

and

$$(3 - \hat{x} \hat{y})/2 - (2^{-M} + 2^{-2-M}) \leq \hat{r} < (3 - \hat{x} \hat{y})/2 - 2^{-M}.$$

□

The following corollary of Lemma 3.5 allows the outputs of *FPU-MUL* to be used as inputs on the next iteration of *FPU-DIV-SQRT*.

**Lemma 3.6.** *Let  $\text{op} \in \{\text{OP-DIV}, \text{OP-SQRT}\}$ ,  $\text{pc} = \text{PC-}^*$ , and  $\text{rc} = \text{RC-NEAR}$ . Assume that  $x$  and  $y$  are normal encodings,  $3/2 < \text{sig}(\hat{x})\text{sig}(\hat{y}) < 3$ , and  $|1 - \hat{x} \hat{y}| < 1/8$ . Then*

- (a) *if  $\text{op} = \text{OP-DIV}$ , then  $3/2 < \text{sig}(\hat{z})\text{sig}(\hat{r}) < 3$ ;*
- (b) *if  $\text{op} = \text{OP-SQRT}$ , then  $3/2 < \text{sig}(\hat{z})\text{sig}(\text{near}(\hat{r}^2, M)) < 3$ .*

*Proof.* Note first that by Theorem 1,  $|1 - \hat{z}| \leq 1/8$ . Now suppose that  $\hat{z} < 1$ . Then  $7/8 \leq \hat{z} < 1$ . If  $\text{op} = \text{OP-DIV}$ , then  $1 \leq \hat{r} < 2 - \hat{x} \hat{y} \leq 9/8$ ; hence  $\text{sig}(\hat{z})\text{sig}(\hat{r}) = 2\hat{z}\hat{r}$  and  $3/2 < 7/4 \leq 2\hat{z}\hat{r} < 9/4 < 3$ . For the case  $\text{op} = \text{OP-SQRT}$ , let  $w = \text{near}(\hat{r}^2, M)$ . Since

$1 \leq \hat{r} < (3 - \hat{x}\hat{y})/2 < 17/16$ ,  $1 \leq \hat{r}^2 < 289/256 < 3/2$ , which implies  $1 \leq w < 3/2$ . Thus,  $\text{sig}(\hat{z})\text{sig}(w) = 2\hat{z}w$  and  $3/2 < 7/4 < 2\hat{z}w < 3$ .

On the other hand, if  $\hat{z} \geq 1$ , then  $1 \leq \hat{z} < 9/8$ . If  $\text{op} = \text{OP-DIV}$ , then  $1 > \hat{r} \geq 2 - \hat{x}\hat{y} - 2^{1-M} \geq 7/8 - 2^{1-M} > 3/4$ , and again  $\text{sig}(\hat{z})\text{sig}(\hat{r}) = 2\hat{z}\hat{r}$ , where  $3/2 < 2\hat{z}\hat{r} < 9/4 < 3$ . If  $\text{op} = \text{OP-SQRT}$ , then  $1 > \hat{r} \geq (3 - \hat{x}\hat{y})/2 - 2^{1-M} \geq 15/16 - 2^{1-M} > 7/8$  and  $1 > \hat{r}^2 > 49/64$ , which implies  $1 > w \geq 49/64 > 3/4$ . Thus,  $\text{sig}(\hat{z})\text{sig}(w) = 2\hat{z}w$  and  $3/2 < 2\hat{z}w \leq 9/4 < 3$ . □

### 3.4. The operation OP-LAST

In the OP-LAST case, the product is rounded to  $\text{mbits}(\text{lastpc}) + 1$  bits, essentially by near rounding.

**Lemma 3.7.** *If  $\text{op} = \text{OP-LAST}$ ,  $\text{pc} = \text{PC-}^*$ ,  $\text{rc} = \text{RC-NEAR}$ ,  $\text{mbits}(\text{lastpc}) = \lambda$ ,  $x$  and  $y$  are normal encodings, and*

$$2^{-2^{17}}(2 - 2^{-\lambda-1}) \leq |\hat{x}\hat{y}| < 2^{2^{17}}(2 - 2^{-\lambda-1}),$$

then

- (a)  $\hat{z}$  is  $(\lambda + 1)$ -exact;
- (b)  $\text{expo}(\hat{x}\hat{y}) \leq \text{expo}(\hat{z})$ ;
- (c)  $|\hat{z} - \hat{x}\hat{y}| \leq 2^{\text{expo}(\hat{x}\hat{y})-\lambda-1}$ .

*Proof.* Note that

$$\text{add} = \text{man-unrounded} + 2^{P-\lambda-2}$$

and by Lemma 2.12,

$$\text{trunc} = 2^{2^M} - 2^{P-\lambda-1} = \text{trunc}'.$$

Let  $\rho = \text{rem}(\text{man-rounded}, 2^P)$ . We shall show that

$$|\rho 2^{\text{round-carryout}} - \text{man-unrounded}| \leq 2^{P-\lambda-2}$$

and that

$$1 - 2^{17} \leq \text{expo}(\hat{x}\hat{y}) + \text{round-carryout} \leq 2^{17},$$

by considering the following two cases.

*Case 1.*  $\text{round-carryout} = 0$

By Lemma 3.3,  $\text{expo}(\text{add}) = \text{expo}(\text{man-rounded}) = P - 1$ ; hence

$$\rho = \text{man-rounded} = \text{add} \ \& \ (2^{2^M} - 2^{P-\lambda-1}) = \text{trunc}(\text{add}, \lambda + 1)$$

by Lemma 2.24. Thus, by Lemma 2.20,

$$\rho \leq \text{add} = \text{man-unrounded} + 2^{P-\lambda-2}$$

and

$$\rho > \text{add} - 2^{(P-1)-(\lambda+1)+1} = \text{man-unrounded} - 2^{P-\lambda-2}.$$

If  $\text{expo}(\hat{x}\hat{y}) = 2^{17}$ , then  $2^{-2^{17}}(2 - 2^{-\lambda-1}) \leq |\hat{x}\hat{y}| < 2^{-2^{17}+1}$ ; hence

$$\text{man-unrounded} = 2^{P-1}\text{sig}(\hat{x}\hat{y}) \geq 2^{P-1}(2 - 2^{-\lambda-1}) = 2^P - 2^{P-\lambda-2},$$

contradicting  $\text{add} < 2^P$ . Thus,  $1 - 2^{17} \leq \text{expo}(\hat{x}\hat{y}) \leq 2^{17}$ .

*Case 2.*  $\text{round-carryout} = 1$

In this case,

$$2^P \leq \text{add} = \text{man-unrounded} + 2^{P-\lambda-2} < 2^P + 2^{P-\lambda-2},$$

which implies

$$|2^P - \text{man-unrounded}| < 2^{P-\lambda-2}$$

as well as

$$\text{rem}(\text{add}, 2^P) < 2^{P-\lambda-2}.$$

Thus, by Lemmas 3.3, 2.9, 2.8, and 2.7,

$$\begin{aligned} \rho &= \text{rem}(2^{P-1} \mid (\text{add} \ \& \ \text{trunc}'), 2^P) = 2^{P-1} \mid (\text{rem}(\text{add}, 2^P) \ \& \ \text{trunc}') \\ &= 2^{P-1} \mid (\text{rem}(\text{add}, 2^P) \ \& \ \text{rem}(\text{trunc}', 2^{P-\lambda-2})) = 2^{P-1} \mid (\text{rem}(\text{add}, 2^P) \ \& \ 0) \\ &= 2^{P-1}, \end{aligned}$$

and therefore

$$|2\rho - \text{man-unrounded}| = |2^P - \text{man-unrounded}| < 2^{P-\lambda-2}.$$

If  $\text{expo}(\hat{x}\hat{y}) = 2^{17}$ , then  $2^{2^{17}} \leq |\hat{x}\hat{y}| < 2^{2^{17}}(2 - 2^{-\lambda-1})$ ; hence

$$\text{man-unrounded} = 2^{P-1} \text{sig}(\hat{x}\hat{y}) < 2^{P-1}(2 - 2^{-\lambda-1}) = 2^P - 2^{P-\lambda-2},$$

contradicting  $\text{add} \geq 2^P$ . Thus,  $1 - 2^{17} \leq \text{expo}(\hat{x}\hat{y}) + 1 \leq 2^{17}$ .

Note that in both cases,  $\rho$  is  $(\lambda + 1)$ -exact; hence so is  $\hat{z}$ , since  $\text{sig}(\hat{z}) = \rho 2^{1-P}$ . Since

$$1 - 2^{17} \leq \text{expo}(\hat{x}\hat{y}) + \text{round-carryout} \leq 2^{17},$$

and  $\text{expo}(\hat{z})$  must lie in the same interval,

$$\text{expo}(\hat{z}) = \text{expo}(\hat{x}\hat{y}) + \text{round-carryout}.$$

Thus,

$$\begin{aligned} |\hat{z} - \hat{x}\hat{y}| &= |\rho 2^{1-P} 2^{\text{expo}(\hat{x}\hat{y}) + \text{round-carryout}} - \text{sig}(\hat{x}\hat{y}) 2^{\text{expo}(\hat{x}\hat{y})}| \\ &= 2^{\text{expo}(\hat{x}\hat{y}) + 1 - P} |\rho 2^{\text{round-carryout}} - \text{man-unrounded}| \\ &\leq 2^{\text{expo}(\hat{x}\hat{y}) + 1 - P} 2^{P-\lambda-2} \\ &= 2^{\text{expo}(\hat{x}\hat{y}) - \lambda - 1}. \end{aligned}$$

□

### 3.5. The operation OP-BACK

In the OP-BACK case, the product is compared, by way of subtraction, to the input  $d$ . The results of the comparison are given by the outputs  $z$  and *inexact*.

**Lemma 3.8.** *If  $\text{op} = \text{OP-BACK}$ ,  $\text{pc} = \text{PC-}^*$ ,  $\text{rc} = \text{RC-CHOP}$ ,  $x$  and  $y$  are normal encodings, and  $|\hat{x}\hat{y} - \hat{d}| < 2^{\text{expo}(\hat{d})-3}$ , then*

- (a)  $|\hat{x}\hat{y}| < |\hat{d}| \Leftrightarrow \text{get-man}(z)[M - 2] = 1$ ;
- (b)  $\hat{x}\hat{y} = \hat{d} \Leftrightarrow \text{get-man}(z)[M - 2 : 0] = \text{inexact} = 0$ .

*Proof.* (a) Since

$$\begin{aligned} \text{rconst-with-overflow} &= \text{comp1}(2^M \text{get-man}(d), 2M) \\ &= 2^{2M} - 2^M \text{get-man}(d) - 1 \end{aligned}$$

and

$$\begin{aligned} \text{rconst-no-overflow} &= \text{shr}(\text{rconst-with-overflow}, 0, 2M) \\ &= \lfloor (2^{2M} - 2^M \text{get-man}(d) - 1) / 2 \rfloor \\ &= 2^{2M-1} - 2^{M-1} \text{get-man}(d) - 1, \end{aligned}$$

we have

$$\text{rconst} = 2^P - 2^{P-M} \text{get-man}(d) - 1,$$

and thus

$$\begin{aligned} \text{add} &= \text{rem}(2^P + \text{man-unrounded} - 2^{P-M} \text{get-man}(d), 2^{P+1}) \\ &= \text{rem}(2^P + 2^{P-1} \text{sig}(\hat{x}\hat{y}) - 2^{P-1} \text{sig}(\hat{d}), 2^{P+1}) \\ &= \text{rem}(2^{P-1}(2 + \text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})), 2^{P+1}) \\ &= 2^{P-1}(2 + \text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})). \end{aligned}$$

Note also that  $\text{trunc}' = \text{trunc} = 2^{2M} - 2^{P-M}$ .

By Lemmas 2.4, 2.5, 2.11, and 3.3,

$$\begin{aligned} \text{get-man}(z)[M - 2 : 0] &= (\text{man-rounded}[P - 1 : P - M])[M - 2 : 0] \\ &= \text{man-rounded}[P - 2 : P - M] \\ &= (\text{add} \ \& \ \text{trunc}')[P - 2 : P - M] \\ &= (2^{P-M} \text{add}[2M - 1 : P - M])[P - 2 : P - M] \\ &= \text{add}[2M - 1 : P - M][M - 2 : 0] \\ &= \text{add}[P - 2 : P - M] \\ &= \rho[P - 2 : P - M], \end{aligned}$$

where  $\rho = \text{rem}(\text{add}, 2^{P-1})$ . In particular, by Lemma 2.5,

$$\begin{aligned} \text{get-man}(z)[M - 2] &= \text{get-man}(z)[M - 2 : 0][M - 2] \\ &= \rho[P - 2 : P - M][M - 2] = \rho[P - 2]. \end{aligned}$$

We must show

$$\rho[P - 2] = 1 \Leftrightarrow |\hat{x}\hat{y}| < |\hat{d}|.$$

Since

$$|\hat{x}\hat{y} - \hat{d}| = |2^{\text{expo}(\hat{x}\hat{y}) - \text{expo}(\hat{d})} \text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})| 2^{\text{expo}(\hat{d})} < 2^{\text{expo}(\hat{d}) - 3},$$

we have

$$|2^{\text{expo}(\hat{x}\hat{y}) - \text{expo}(\hat{d})} \text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})| < 2^{-3},$$

which implies  $|\text{expo}(\hat{x}\hat{y}) - \text{expo}(\hat{d})| \leq 1$ . Thus, we have three cases to consider.

*Case 1.*  $\text{expo}(\hat{x}\hat{y}) = \text{expo}(\hat{d})$

In this case,  $|\text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})| < 2^{-3}$ .

Suppose first that  $|\hat{x}\hat{y}| < |\hat{d}|$ . Then  $\text{sig}(\hat{x}\hat{y}) < \text{sig}(\hat{d})$  and

$$2^P > \text{add} = 2^{P-1}(2 + \text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})) > 2^{P-1}(2 - 2^{-3}) > 2^{P-1} + 2^{P-2}.$$

Thus,

$$2^{P-2} < \rho < 2^{P-1},$$

and  $\rho[P - 2] = 1$  by Lemma 2.2.

On the other hand, if  $|\hat{x}\hat{y}| \geq |\hat{d}|$ , then  $\text{sig}(\hat{x}\hat{y}) \geq \text{sig}(\hat{d})$  and

$$2^P \leq \text{add} < 2^{P-1}(2 + 2^{-3}) < 2^P + 2^{P-2},$$

hence  $\rho < 2^{P-2}$  and  $\rho[P - 2] = 0$ .

Case 2.  $\text{expo}(\hat{x}\hat{y}) = \text{expo}(\hat{d}) + 1$

Here,  $|\hat{x}\hat{y}| > |\hat{d}|$  and

$$0 < 2\text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d}) < 2^{-3}.$$

Thus,

$$\text{sig}(\hat{x}\hat{y}) < \frac{1}{2}\text{sig}(\hat{d}) + 2^{-4} \leq 1 + 2^{-4}$$

and

$$\text{sig}(\hat{d}) > 2\text{sig}(\hat{x}\hat{y}) - 2^{-3} \geq 2 - 2^{-3}.$$

It follows that

$$\text{add} < 2^{P-1}(2 + 1 + 2^{-4} - 2 + 2^{-3}) < 2^{P-1} + 2^{P-2}.$$

But  $\text{add} > 2^{P-1}(2 + 1 - 2) = 2^{P-1}$ ; hence  $\rho < 2^{P-2}$  and  $\rho[P - 2] = 0$ .

Case 3.  $\text{expo}(\hat{x}\hat{y}) = \text{expo}(\hat{d}) - 1$

In this case,  $|\hat{x}\hat{y}| < |\hat{d}|$  and

$$0 < \text{sig}(\hat{d}) - \frac{1}{2}\text{sig}(\hat{x}\hat{y}) < 2^{-3}.$$

Thus,  $\text{sig}(\hat{d}) < 1 + 2^{-3}$ ,  $\text{sig}(\hat{x}\hat{y}) > 2 - 2^{-2}$ , and

$$\text{add} > 2^{P-1}(2 + 2 - 2^{-2} - 1 - 2^{-3}) > 3 \cdot 2^{P-1} - 2^{P-2} = 2 \cdot 2^{P-1} + 2^{P-2}.$$

But

$$\text{add} < 2^{P-1}(2 + 2 - 1) = 3 \cdot 2^{P-1};$$

hence  $\rho > 2^{P-2}$  and  $\rho[P - 2] = 1$ .

(b) Note that by Lemmas 3.1 and 3.2,  $\text{inexact} = 0$  iff  $\hat{x}\hat{y}$  is  $M$ -exact. Thus, if  $\hat{x}\hat{y} = \hat{d}$ , then  $\text{inexact} = 0$  and  $\text{add} = 2^P$ , which implies  $\rho = 0$ , and hence  $\text{get-man}(z)[M - 2 : 0] = 0$ .

Conversely, suppose

$$\text{get-man}(z)[M - 2 : 0] = \rho[P - 2 : P - M] = \text{inexact} = 0.$$

Then  $\text{sig}(\hat{x}\hat{y})$  is  $M$ -exact, i.e.,  $2^{M-1}\text{sig}(\hat{x}\hat{y}) \in \mathbb{Z}$ , hence  $2^{P-1}\text{sig}(\hat{x}\hat{y})$  is divisible by  $2^{P-M}$ . Similarly,  $2^{P-1}\text{sig}(\hat{d})$  is divisible by  $2^{P-M}$ , and hence, so are  $\text{add}$  and  $\rho$ . Thus,

$$\rho = (\rho/2^{P-M})2^{P-M} = \lfloor \rho/2^{P-M} \rfloor 2^{P-M} = \rho[P - 2 : P - M]2^{P-M} = 0.$$

Since  $\hat{x}\hat{y} = -\hat{d}$  is impossible, we need only show  $|\hat{x}\hat{y}| = |\hat{d}|$ . In view of (a), we may assume  $|\hat{x}\hat{y}| \geq |\hat{d}|$ . Thus, there are two cases to consider.

Case 1.  $\text{expo}(\hat{x}\hat{y}) = \text{expo}(\hat{d})$

In this case,  $\text{sig}(\hat{x}\hat{y}) \geq \text{sig}(\hat{d})$ , which implies

$$\rho = 2^{P-1}(\text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})) = 0,$$

hence  $\text{sig}(\hat{x}\hat{y}) = \text{sig}(\hat{d})$  and  $|\hat{x}\hat{y}| = |\hat{d}|$ .

Case 2.  $\text{expo}(\hat{x}\hat{y}) = \text{expo}(\hat{d}) + 1$

If this were to occur, then we would have

$$\rho = 2^{P-1}(1 + \text{sig}(\hat{x}\hat{y}) - \text{sig}(\hat{d})) = 0,$$

implying  $\text{sig}(\hat{d}) = 1 + \text{sig}(\hat{x}\hat{y}) \geq 2$ , which is impossible. □

#### 4. Division and square root

##### 4.1. The program FPU-DIV-SQRT

The hardware for division and square root is represented by the program *FPU-DIV-SQRT*, shown in Figures 3 and 4. Our analysis will be based on an execution of

$$FPU-DIV-SQRT(op, pc, rc, a, b, z),$$

with inputs as follows.

- (a)  $op \in \{\text{OP-DIV}, \text{OP-SQRT}\}$ ;
- (b)  $pc$  is an external precision control specifier;
- (c)  $rc$  is a rounding control specifier;
- (d)  $a$  and  $b$  are normal encodings.

In the case  $op = \text{OP-DIV}$ , the output  $z$  represents an appropriately rounded approximation of the quotient  $\hat{a}/\hat{b}$ ; when  $op = \text{OP-SQRT}$ ,  $a$  is ignored and an approximation of  $\sqrt{\hat{b}}$  is returned.

Both operations are based on Goldschmidt's Algorithm [1], a variant of Newton-Raphson approximation. Our analysis of division will involve a sequence  $\xi_0, \xi_1, \xi_2, \xi_3$  of approximations to  $1/\hat{b}$ , where  $\xi_0$  is derived from a table and the other  $\xi_i$  are computed by three successive Newton-Raphson iterations. The square root involves a similar sequence of approximations to  $1/\sqrt{\hat{b}}$ .

Although the algorithm does not explicitly compute the  $\xi_i$  for  $i > 0$ , a sequence of calls to *FPU-MUL* produces an encoding  $q$  of either  $\hat{a}\xi_i$  or  $\hat{b}\xi_i$ , modulo rounding error, according to whether  $op = \text{OP-DIV}$  or  $op = \text{OP-SQRT}$ , where (a)  $i = 1$  if  $pc = \text{PC-32}$ , (b)  $i = 2$  if  $pc = \text{PC-64}$ , and (c)  $i = 3$  if  $pc = \text{PC-80}$  or  $pc = \text{PC-87}$ . Lemmas 4.9 and 4.13 give estimates of the errors  $|\hat{q} - \hat{a}/\hat{b}|$  and  $|\hat{q} - \sqrt{\hat{b}}|$ . Note that the constraint  $M \geq 75$  on the multiplier width is required in the proofs of these lemmas.

The approximation  $\hat{q}$  is compared to the exact value by means of a final call to *FPU-MUL* with  $op = \text{OP-BACK}$ . Using the results of this comparison,  $q$  is then adjusted to produce the correctly rounded result  $z$ . The correctness of this result is guaranteed by Theorems 2 and 3.

##### 4.2. Initial approximation

The initial approximation  $x_0$  to the reciprocal of  $b$ , in the case  $op = \text{OP-DIV}$ , is derived from a pair of tables, each consisting of  $2^{10}$  bit vectors, which we represent by the functions *recip-rom-p* and *recip-rom-n*. If  $\text{sig}(\hat{b})$  has the binary representation  $1.b_1b_2b_3 \dots$ , then the bit vectors

$$b_1b_2 \dots b_9b_{10} = \text{get-man}(b)[M - 2 : M - 11]$$

and

$$b_1 \dots b_5b_{11} \dots b_{15} = \text{cat}(\text{get-man}(b)[M - 2 : M - 6], \text{get-man}(b)[M - 12 : M - 16], 5)$$



**Program** *FPU-DIV-SQRT*(*op,pc,rc,a,b,z*):

if *op* = OP-DIV then

{*sign* ← get-sign(*a*) ^ get-sign(*b*);

*p*-value ← recip-rom-*p*(get-man(*b*)[*M* - 2 : *M* - 11]);

*n*-value ← recip-rom-*n*(cat(get-man(*b*)[*M* - 2 : *M* - 6],  
get-man(*b*)[*M* - 12 : *M* - 16],  
5));

estimate ← (*p*-value + *n*-value)[16 : 0];

*x*<sub>0</sub> ← (get-sign(*b*),

$2^{M-17}$ estimate |  $2^{M-1}$ ,

( $2^{18} - 2 + \text{comp1}(\text{get-expo}(b), 18) + \text{estimate}[16]$ )[17 : 0]);

*FPU-MUL*(OP-DIV, PC-\*, NIL, RC-NEAR, *b*, *x*<sub>0</sub>, *d*<sub>0</sub>, *r*<sub>0</sub>, NIL, NIL);

*FPU-MUL*(OP-MUL, PC-\*, NIL, RC-NEAR, *a*, *x*<sub>0</sub>, *n*<sub>0</sub>, NIL, NIL, NIL);

if *pc* = PC-32

then *FPU-MUL*(OP-LAST, PC-\*, *pc*, RC-NEAR, *n*<sub>0</sub>, *r*<sub>0</sub>, *q*, NIL, NIL, NIL)

else {*FPU-MUL*(OP-DIV, PC-\*, NIL, RC-NEAR, *d*<sub>0</sub>, *r*<sub>0</sub>, *d*<sub>1</sub>, *r*<sub>1</sub>, NIL, NIL);

*FPU-MUL*(OP-MUL, PC-\*, NIL, RC-NEAR, *n*<sub>0</sub>, *r*<sub>0</sub>, *n*<sub>1</sub>, NIL, NIL, NIL);

if *pc* = PC-64

then *FPU-MUL*(OP-LAST, PC-\*, *pc*, RC-NEAR, *n*<sub>1</sub>, *r*<sub>1</sub>, *q*, NIL, NIL, NIL)

else {*FPU-MUL*(OP-DIV, PC-\*, NIL, RC-NEAR, *d*<sub>1</sub>, *r*<sub>1</sub>, *d*<sub>2</sub>, *r*<sub>2</sub>, NIL, NIL);

*FPU-MUL*(OP-MUL, PC-\*, NIL, RC-NEAR, *n*<sub>1</sub>, *r*<sub>1</sub>, *n*<sub>2</sub>, NIL, NIL, NIL);

*FPU-MUL*(OP-LAST, PC-\*, *pc*, RC-NEAR, *n*<sub>2</sub>, *r*<sub>2</sub>, *q*, NIL, NIL, NIL)};

*FPU-MUL*(OP-BACK, PC-\*, NIL, RC-CHOP, *b*, *q*, rem, NIL, *a*, inexact)

else if *op* = OP-DIV-SQRT then

{*sign* ← 0;

*p*-value ← sqrt-rom-*p*(cat(get-expo(*b*)[0], get-man(*b*)[*M* - 2 : *M* - 11], 10));

*n*-value ← sqrt-rom-*n*(cat(get-expo(*b*)[0],

cat(get-man(*b*)[*M* - 2 : *M* - 6],

get-man(*b*)[*M* - 12 : *M* - 16],

5),

10));

estimate ← (*p*-value + *n*-value)[16 : 0];

*x*<sub>0</sub> ← (get-sign(*b*),

$2^{M-17}$ estimate |  $2^{M-1}$ ,

shr(( $2^{18} + 2^{17} - 3 + \text{comp1}(\text{get-expo}(b), 19) + \text{estimate}[16]$ )[18 : 0], 0, 19));

*FPU-MUL*(OP-MUL, PC-\*, NIL, RC-NEAR, *x*<sub>0</sub>, *x*<sub>0</sub>, *t*<sub>0</sub>, NIL, NIL, NIL);

*FPU-MUL*(OP-MUL, PC-\*, NIL, RC-NEAR, *b*, *x*<sub>0</sub>, *d*<sub>0</sub>, NIL, NIL, NIL);

*FPU-MUL*(OP-SQRT, PC-\*, NIL, RC-NEAR, *b*, *t*<sub>0</sub>, *n*<sub>0</sub>, *r*<sub>0</sub>, NIL, NIL);

Figure 3: *FPU-DIV-SQRT*

```

if pc = PC-32
  then  $FPU-MUL(OP-LAST, PC-*, pc, RC-NEAR, d_0, r_0, q, NIL, NIL, NIL)$ 
  else { $FPU-MUL(OP-MUL, PC-*, NIL, RC-NEAR, r_0, r_0, t_1, NIL, NIL, NIL)$ ;
         $FPU-MUL(OP-MUL, PC-*, NIL, RC-NEAR, d_0, r_0, d_1, NIL, NIL, NIL)$ ;
         $FPU-MUL(OP-SQRT, PC-*, NIL, RC-NEAR, n_0, t_1, n_1, r_1, NIL, NIL)$ ;
        if pc = PC-64
          then  $FPU-MUL(OP-LAST, PC-*, pc, RC-NEAR, d_1, r_1, q, NIL, NIL, NIL)$ 
          else { $FPU-MUL(OP-MUL, PC-*, NIL, RC-NEAR, r_1, r_1, t_2, NIL, NIL, NIL)$ ;
                 $FPU-MUL(OP-MUL, PC-*, NIL, RC-NEAR, d_1, r_1, d_2, NIL, NIL, NIL)$ ;
                 $FPU-MUL(OP-SQRT, PC-*, NIL, RC-NEAR, n_1, t_2, n_2, r_2, NIL, NIL)$ ;
                 $FPU-MUL(OP-LAST, PC-*, pc, RC-NEAR, d_2, r_2, q, NIL, NIL, NIL)$ };
         $FPU-MUL(OP-BACK, PC-*, NIL, RC-CHOP, q, q, rem, NIL, b, inexact)$ };

if get-man(rem)[ $M - 2 : 0$ ] = 0
  then rem-zero  $\leftarrow$  compl(inexact, 1)
  else rem-zero  $\leftarrow$  0;
rem-neg  $\leftarrow$  compl(get-man(rem)[ $M - 2$ ], 1) & compl(rem-zero, 1);
rem-pos  $\leftarrow$  get-man(rem)[ $M - 2$ ];
q-lsb  $\leftarrow$  get-man(q)[ $M - \text{mbits}(pc)$ ];
q-guard  $\leftarrow$  get-man(q)[ $M - \text{mbits}(pc) - 1$ ];
if op = OP-DIV  $\wedge$  get-man(a) = 0 then
  z  $\leftarrow$  (sign, 0, get-expo(a))
else if op = OP-SQRT  $\wedge$  get-man(b) = 0 then
  z  $\leftarrow$  (sign, 0, get-expo(b))
else if ((rc = RC-POS  $\wedge$  sign = 1)  $\vee$  (rc = RC-NEG  $\wedge$  sign = 0)  $\vee$  rc = RC-CHOP)
   $\wedge$  q-guard = 0  $\wedge$  rem-neg = 1 then
  if get-man(q) & ( $2^M - 2^{M-\text{mbits}(pc)}$ ) =  $2^{M-1}$ 
    then z  $\leftarrow$  (sign,  $2^M - 2^{M-\text{mbits}(pc)}$ , dec1(get-expo(q), 18))
    else z  $\leftarrow$  (sign,
      ((get-man(q) & ( $2^M - 2^{M-\text{mbits}(pc)}$ ))) +  $2^M - 2^{M-\text{mbits}(pc)}$ )[ $M - 1 : 0$ ],
      get-expo(q))
else if (((rc = RC-POS  $\wedge$  sign = 0)  $\vee$  (rc = RC-NEG  $\wedge$  sign = 1))
   $\wedge$  (q-guard = 1  $\vee$  rem-pos = 1))
   $\vee$  (rc = RC-NEAR  $\wedge$  q-guard = 1  $\wedge$  rem-pos = 1)
   $\vee$  (rc = RC-NEAR  $\wedge$  q-guard = 1  $\wedge$  rem-zero = 1  $\wedge$  q-lsb = 1) then
  if get-man(q) & ( $2^M - 2^{M-\text{mbits}(pc)}$ ) =  $2^M - 2^{M-\text{mbits}(pc)}$ 
    then z  $\leftarrow$  (sign,  $2^{M-1}$ , (get-expo(q) + 1)[17 : 0])
    else z  $\leftarrow$  (sign,
      ((get-man(q) & ( $2^M - 2^{M-\text{mbits}(pc)}$ ))) +  $2^{M-\text{mbits}(pc)}$ )[ $M - 1 : 0$ ],
      get-expo(q))
else z  $\leftarrow$  (sign, get-man(q) & ( $2^M - 2^{M-\text{mbits}(pc)}$ ), get-expo(q)).

```

Figure 4:  $FPU-DIV-SQRT$  (continued)

are used as indices into these tables. The results are added and the 16-bit sum is appended to a leading 1 and  $M - 17$  trailing 0's to produce  $\text{get-man}(x_0)$ . For  $\text{op} = \text{OP-SQRT}$ , a separate pair of tables, represented by the functions  $\text{sqrt-rom-p}$  and  $\text{sqrt-rom-n}$ , is similarly used to derive an initial approximation to the reciprocal of the square root of  $b$ .

The functions  $R_0$ ,  $S_0$ , and  $S_1$ , which are defined in terms of these functions, represent the computation of  $\text{get-man}(x_0)$  in the three cases listed in Lemma 4.4 below.

**Definition 4.1.** For all  $i \in \mathbb{N}$ ,

- (a)  $R_0(i) = 2^{16} + \text{recip-rom-p}(i[14 : 5]) + \text{recip-rom-n}(\text{cat}(i[14 : 10], i[4 : 0], 5))$ ;
- (b)  $S_0(i) = 2^{16} + \text{sqrt-rom-p}(i[14 : 5]) + \text{sqrt-rom-n}(\text{cat}(i[14 : 10], i[4 : 0], 5))$ ;
- (c)  $S_1(i) = 2^{16} + \text{sqrt-rom-p}(2^{10} + i[14 : 5]) + \text{sqrt-rom-n}(2^{10} + \text{cat}(i[14 : 10], i[4 : 0], 5))$ .

While space does not allow a complete listing of the tables here, we list instead the following three lemmas, which contain all required relevant information, and which have all been verified by direct computation, using ACL2.

**Lemma 4.1.** For all  $i \in \mathbb{N}$ , if  $i < 2^{15}$ , then  $R_0(i) \in \mathbb{N}$ ,  $S_0(i) \in \mathbb{N}$ ,  $S_1(i) \in \mathbb{N}$ , and

$$\text{expo}(R_0(i)) = \text{expo}(S_0(i)) = \text{expo}(S_1(i)) = 16.$$

**Lemma 4.2.** For all  $i \in \mathbb{N}$ , if  $i < 2^{15}$ , then

- (a)  $2^{32} - 3 \cdot 2^{16} < R_0(i)(2^{15} + i) < R_0(i)(2^{15} + i + 1) < 2^{32} + 3 \cdot 2^{16}$ ;
- (b)  $2^{48} - 3 \cdot 2^{32} < S_0(i)^2(2^{15} + i) < S_0(i)^2(2^{15} + i + 1) < 2^{48} + 3 \cdot 2^{32}$ ;
- (c)  $2^{49} - 3 \cdot 2^{33} < S_1(i)^2(2^{15} + i) < S_1(i)^2(2^{15} + i + 1) < 2^{49} + 3 \cdot 2^{33}$ .

**Lemma 4.3.** For all  $i \in \mathbb{N}$ , if  $i < 2^{15}$ , then  $S_0(i)^2 < 2^{33} \leq S_1(i)^2$ .

The relationship between  $x_0$  and  $b$  may be described in terms of  $R_0$ ,  $S_0$ , and  $S_1$ .

**Lemma 4.4.** Let  $I = \text{get-man}(b)[M - 2 : M - 16]$ . Assume that if  $\text{op} = \text{OP-DIV}$ , then  $\text{get-expo}(b) \leq 2^{18} - 3$ . Then  $x_0$  is normal and

- (a)  $\text{sgn}(\hat{x}_0) = \begin{cases} \text{sgn}(\hat{b}) & \text{if op} = \text{OP-DIV} \\ 1 & \text{if op} = \text{OP-SQRT}; \end{cases}$
- (b)  $\text{sig}(\hat{x}_0) = \begin{cases} 2^{-16}R_0(I) & \text{if op} = \text{OP-DIV} \\ 2^{-16}S_0(I) & \text{if op} = \text{OP-SQRT and get-expo}(b)[0] = 0 \\ 2^{-16}S_1(I) & \text{if op} = \text{OP-SQRT and get-expo}(b)[0] = 1; \end{cases}$
- (c)  $\text{expo}(\hat{x}_0) = \begin{cases} -\text{expo}(\hat{b}) - 1 & \text{if op} = \text{OP-DIV} \\ -\lfloor \text{expo}(\hat{b})/2 \rfloor - 1 & \text{if op} = \text{OP-SQRT}. \end{cases}$

*Proof.* First consider the case  $\text{op} = \text{OP-DIV}$ . By Lemma 2.5,

$$\text{get-man}(b)[M - 2 : M - 11] = \text{get-man}(b)[M - 2 : M - 16][14 : 5] = I[14 : 5],$$

hence  $p\text{-value} = \text{recip-rom-p}(I[14 : 5])$ . Similarly,

$$n\text{-value} = \text{recip-rom-n}(\text{cat}(I[14 : 10], I[4 : 0], 5)).$$

By Lemma 4.1,

$$p\text{-value} + n\text{-value} = R_0(I) - 2^{16} < 2^{17} - 2^{16} = 2^{16},$$

hence

$$\text{estimate} = p\text{-value} + n\text{-value} < 2^{16}$$

and by Lemma 2.8,

$$\begin{aligned} \text{get-man}(x_0) &= 2^{M-17} \text{estimate} \mid 2^{M-1} = 2^{M-17} (\text{estimate} \mid 2^{16}) \\ &= 2^{M-17} (\text{estimate} + 2^{16}) = 2^{M-17} R_0(I). \end{aligned}$$

Since  $\text{estimate}[16] = 0$  and  $\text{get-expo}(b) \leq 2^{18} - 3$ ,

$$\text{get-expo}(x_0) = \text{rem}(2^{18} - 2 + 2^{18} - \text{get-expo}(b) - 1, 2^{18}) = 2^{18} - 3 - \text{get-expo}(b).$$

The OP-DIV case now follows easily from Lemmas 4.1 and 2.15.

In the case  $\text{op} = \text{OP-SQRT}$ , we may similarly show that  $\text{get-man}(x_0) = 2^{M-17} S_j(I)$ , where  $j = \text{get-expo}(b)[0]$ . Now

$$\begin{aligned} &(2^{18} + 2^{17} - 3 + \text{comp1}(\text{get-expo}(b), 19) + \text{estimate}[16])[18 : 0] \\ &= (2^{18} + 2^{17} - 3 + \text{comp1}(\text{get-expo}(b), 19))[18 : 0] \\ &= \text{rem}(2^{18} + 2^{17} - 3 + \text{comp1}(\text{get-expo}(b), 19), 2^{19}) \\ &= \text{rem}(2^{18} + 2^{17} - 3 + 2^{19} - \text{get-expo}(b) - 1, 2^{19}) \\ &= \text{rem}(2^{18} + 2^{17} - 3 + 2^{19} - (\text{expo}(\hat{b}) + 2^{17} - 1) - 1, 2^{19}) \\ &= \text{rem}(2^{18} - \text{expo}(\hat{b}) - 3, 2^{19}) \\ &= 2^{18} - \text{expo}(\hat{b}) - 3. \end{aligned}$$

Thus,

$$\begin{aligned} \text{get-expo}(x_0) &= \text{shr}(2^{18} - \text{expo}(\hat{b}) - 3, 0, 19) \\ &= \lfloor (2^{18} - \text{expo}(\hat{b}) - 3)/2 \rfloor \\ &= 2^{17} - 1 + \lfloor -(\text{expo}(\hat{b}) + 1)/2 \rfloor, \end{aligned}$$

and

$$\text{expo}(\hat{x}_0) = \lfloor -(\text{expo}(\hat{b}) + 1)/2 \rfloor = -\lfloor \text{expo}(\hat{b})/2 \rfloor - 1.$$

□

The error associated with  $x_0$  is characterized by the next two lemmas, which also establish the bounds required by Lemma 3.5.

**Lemma 4.5.** *If  $\text{op} = \text{OP-DIV}$  and  $\text{get-expo}(b) \leq 2^{18} - 3$ , then*

$$(a) |1 - \hat{x}_0 \hat{b}| < 3 \cdot 2^{-16}; \quad (b) 3/2 < \text{sig}(\hat{x}_0) \text{sig}(\hat{b}) < 3.$$

*Proof.* (a) By Lemma 4.4,

$$\hat{x}_0 \hat{b} = \text{sig}(\hat{x}_0) \text{sig}(\hat{b}) 2^{\text{expo}(\hat{x}_0) + \text{expo}(\hat{b})} = \text{sig}(\hat{x}_0) \text{sig}(\hat{b})/2.$$

Let  $I = \text{get-man}(b)[M - 2 : M - 16]$ . Since  $2^{M-1} \leq \text{get-man}(b) < 2^M$ ,

$$\begin{aligned} I &= \lfloor \text{rem}(\text{get-man}(b), 2^{M-1})/2^{M-16} \rfloor = \lfloor (\text{get-man}(b) - 2^{M-1})/2^{M-16} \rfloor \\ &= \lfloor \text{get-man}(b)/2^{M-16} - 2^{15} \rfloor, \end{aligned}$$

hence

$$\text{get-man}(b)/2^{M-16} - 2^{15} - 1 < I \leq \text{get-man}(b)/2^{M-16} - 2^{15},$$

which along with Lemma 2.15, implies

$$2^{-15}(2^{15} + I) \leq \text{sig}(\hat{b}) < 2^{-15}(2^{15} + I + 1).$$

Thus, by Lemmas 4.4 and 4.2,

$$1 - 3 \cdot 2^{-16} < 2^{-32}R_0(I)(2^{15} + I) \leq \hat{x}_0\hat{b} < 2^{-32}R_0(I)(2^{15} + I + 1) < 1 + 3 \cdot 2^{-16}.$$

(b) This follows from (a) and the observation that  $\text{sig}(\hat{x}_0)\text{sig}(\hat{b}) = 2\hat{x}_0\hat{b}$ . □

**Lemma 4.6.** *If  $\text{op} = \text{OP-SQRT}$ ,  $\hat{b} > 0$ , and  $\text{get-expo}(b) \leq 2^{18} - 3$ , then*

- (a)  $|1 - \hat{x}_0^2\hat{b}| < 3 \cdot 2^{-16}$ ;
- (b)  $3/2 < \text{sig}(\hat{x}_0^2)\text{sig}(\hat{b}) < 3$ ;
- (c)  $\hat{x}_0^2$  is representable.

*Proof.* Let  $I = \text{get-man}(b)[M - 2 : M - 16]$  and  $\text{expo}(\hat{b}) = 2r + s$ , where  $0 \leq s \leq 1$ .  
*Case 1.  $s = 0$*

(a) In this case,  $\text{get-expo}(b)[0] = 1$ . By Lemma 4.4,

$$\begin{aligned} \hat{x}_0^2\hat{b} &= \text{sig}(\hat{x}_0)^2\text{sig}(\hat{b})2^{2\text{expo}(\hat{x}_0)+\text{expo}(\hat{b})} = \text{sig}(\hat{x}_0)^2\text{sig}(\hat{b})2^{2(-r-1)+2r} \\ &= \text{sig}(\hat{x}_0)^2\text{sig}(\hat{b})/4 = 2^{-34}S_1(I)^2\text{sig}(\hat{b}). \end{aligned}$$

Thus, by Lemma 4.2,

$$1 - 3 \cdot 2^{-16} < 2^{-49}S_1(I)^2(2^{15} + I) \leq \hat{x}_0^2\hat{b} < 2^{-49}S_1(I)^2(2^{15} + I + 1) < 1 + 3 \cdot 2^{-16}.$$

(b) By Lemmas 4.4 and 4.3,  $\text{sig}(\hat{x}_0)^2 = 2^{-32}S_1(I)^2 \geq 2$ , which implies  $\text{sig}(\hat{x}_0^2) = \text{sig}(\hat{x}_0)^2/2$ . Thus,

$$\hat{x}_0^2\hat{b} = \text{sig}(\hat{x}_0)^2\text{sig}(\hat{b})/4 = \text{sig}(\hat{x}_0^2)\text{sig}(\hat{b})/2.$$

The claim now follows from (a).

(c) By Lemmas 4.1 and 4.4,  $\hat{x}_0$  is 17-exact, and it follows that  $\hat{x}_0^2$  is  $M$ -exact. Since  $\text{expo}(\hat{b}) \geq 1 - 2^{17}$ ,

$$\text{expo}(\hat{x}_0) \leq -\lfloor(1 - 2^{17})/2\rfloor - 1 = 2^{16} - 1$$

and

$$\text{expo}(\hat{x}_0^2) \leq 2\text{expo}(\hat{x}_0) + 1 \leq 2^{17} - 1.$$

But since  $\text{expo}(\hat{b}) = \text{get-expo}(b) - (2^{17} - 1) \leq (2^{18} - 3) - (2^{17} - 1) = 2^{17} - 2$ ,

$$\hat{x}_0^2 = \text{sig}(\hat{x}_0^2)\text{sig}(\hat{b})/2\hat{b} \geq \text{sig}(\hat{b})/2\hat{b} = 2^{-1-\text{expo}(\hat{b})} \geq 2^{1-2^{17}},$$

hence  $\text{expo}(\hat{x}_0^2) \geq 1 - 2^{17}$ .

*Case 2.  $s = 1$*

(a) In this case,  $\text{get-expo}(b)[0] = 0$ . By Lemma 4.4,

$$\begin{aligned} \hat{x}_0^2\hat{b} &= \text{sig}(\hat{x}_0)^2\text{sig}(\hat{b})2^{2\text{expo}(\hat{x}_0)+\text{expo}(\hat{b})} = \text{sig}(\hat{x}_0)^2\text{sig}(\hat{b})2^{2(-r-1)+2r+1} \\ &= \text{sig}(\hat{x}_0)^2\text{sig}(\hat{b})/2 = 2^{-33}S_0(I)^2\text{sig}(\hat{b}). \end{aligned}$$

Thus, by Lemma 4.2,

$$1 - 3 \cdot 2^{-16} < 2^{-48}S_0(I)^2(2^{15} + I) \leq \hat{x}_0^2\hat{b} < 2^{-48}S_0(I)^2(2^{15} + I + 1) < 1 + 3 \cdot 2^{-16}.$$

(b) By Lemmas 4.4 and 4.3,  $\text{sig}(\hat{x}_0)^2 = 2^{-32} S_0(I)^2 < 2$ , which implies  $\text{sig}(\hat{x}_0^2) = \text{sig}(\hat{x}_0)^2$ . Thus,

$$\hat{x}_0^2 \hat{b} = \text{sig}(\hat{x}_0)^2 \text{sig}(\hat{b})/2 = \text{sig}(\hat{x}_0^2) \text{sig}(\hat{b})/2.$$

(c) As in Case 1,  $\hat{x}_0^2$  is  $M$ -exact and  $\text{expo}(\hat{x}_0^2) \leq 2^{17} - 1$ . Since  $\text{expo}(\hat{b}) \leq 2^{17} - 2$  and  $\text{expo}(\hat{b})$  is odd,  $\text{expo}(\hat{b}) \leq 2^{17} - 3$ ; hence

$$\text{expo}(\hat{x}_0) \geq -\lfloor (2^{17} - 3)/2 \rfloor - 1 = 1 - 2^{16}$$

and

$$\text{expo}(\hat{x}_0^2) \geq 2\text{expo}(\hat{x}_0) \geq 2 - 2^{17}.$$

□

### 4.3. The operation OP-DIV

Given an initial approximation  $\xi_0$  of  $1/\hat{b}$ , the Newton-Raphson formula

$$\xi_i = \xi_{i-1}(2 - \hat{b}\xi_{i-1})$$

gives a converging sequence of approximations  $\xi_1, \xi_2, \dots$ . The relative error of  $\xi_i$  is

$$\left| \frac{1/\hat{b} - \xi_i}{1/\hat{b}} \right| = |1 - \hat{b}\xi_i|.$$

Thus, the following lemma (which is proved by simple arithmetic) shows that this sequence is quadratically convergent.

**Lemma 4.7.** *Let  $b, x \in \mathbb{Q}$  and let  $y = x(2 - bx)$ . Then  $1 - by = (1 - bx)^2$ .*

Using Lemma 4.7, we shall derive an error estimate for  $\hat{q}$  as an approximation of  $\hat{a}/\hat{b}$ . First, we prove the following technical lemma.

**Lemma 4.8.** *Assume  $\hat{q}$  is  $(\mu + 1)$ -exact, where  $\mu \geq 1$ , and  $\hat{q} \neq 0$ . Let  $\zeta \in \mathbb{Q}$  satisfy*

$$\text{expo}(\zeta) \leq \text{expo}(\hat{q}),$$

$$|\hat{q} - \zeta| \leq 2^{\text{expo}(\zeta) - \mu - 1},$$

and

$$|\hat{a}/\hat{b} - \zeta| < 2^{\text{expo}(\hat{a}/\hat{b}) - \mu - 2}.$$

Then

$$|\hat{q} - \hat{a}/\hat{b}| < 2^{\min(\text{expo}(\hat{q}), \text{expo}(\hat{a}/\hat{b})) - \mu}.$$

*Proof.* First note that  $|\hat{q}| \geq \frac{3}{4}|\zeta| > \frac{9}{16}|\hat{a}/\hat{b}|$ , hence  $\text{expo}(\hat{q}) \geq \text{expo}(\hat{a}/\hat{b}) - 1$ . Since

$$|\hat{q} - \hat{a}/\hat{b}| \leq |\hat{q} - \zeta| + |\hat{a}/\hat{b} - \zeta| < 2^{\text{expo}(\hat{q}) - \mu - 1} + 2^{\text{expo}(\hat{a}/\hat{b}) - \mu - 2} = 2^{\text{expo}(\hat{q}) - \mu},$$

we may assume  $\text{expo}(\hat{a}/\hat{b}) < \text{expo}(\hat{q})$ . But  $|\hat{a}/\hat{b}| > |\hat{q}|/2$ , hence  $\text{expo}(\hat{a}/\hat{b}) = \text{expo}(\hat{q}) - 1$ .

We may also assume  $\text{expo}(\zeta) = \text{expo}(\hat{q})$ , for otherwise  $\text{expo}(\zeta) \leq \text{expo}(\hat{a}/\hat{b})$  and

$$|\hat{q} - \hat{a}/\hat{b}| \leq |\hat{q} - \zeta| + |\hat{a}/\hat{b} - \zeta| < 2^{\text{expo}(\zeta) - \mu - 1} + 2^{\text{expo}(\hat{a}/\hat{b}) - \mu - 2} \leq 2^{\text{expo}(\hat{a}/\hat{b}) - \mu}.$$

If  $|\hat{q}| > 2^{\text{expo}(\hat{q})}$ , then  $|\hat{q}| \geq 2^{\text{expo}(\hat{q})} + 2^{\text{expo}(\hat{q})-\mu}$  by Lemma 2.13, and

$$|\hat{q} - \hat{a}/\hat{b}| \geq |\hat{q}| - |\hat{a}/\hat{b}| > 2^{\text{expo}(\hat{q})} + 2^{\text{expo}(\hat{q})-\mu} - 2^{\text{expo}(\hat{a}/\hat{b})+1} = 2^{\text{expo}(\hat{q})-\mu}.$$

Therefore,  $|\hat{q}| = 2^{\text{expo}(\hat{q})}$ , which implies  $|\zeta| \geq |\hat{q}|$  and

$$|q - \hat{a}/\hat{b}| = |\hat{q}| - |\hat{a}/\hat{b}| \leq |\zeta| - |\hat{a}/\hat{b}| \leq |\zeta - \hat{a}/\hat{b}| < 2^{\text{expo}(\hat{a}/\hat{b})-\mu}.$$

□

We shall assume here that  $\hat{a}$  and  $\hat{b}$  are both positive; this assumption will be relieved as in the proof of Theorem 1 .

**Lemma 4.9.** Assume  $\text{op} = \text{OP-DIV}$ ,  $\hat{a} > 0$ ,  $\hat{b} > 0$ ,  $\text{expo}(\hat{b}) \leq 2^{17} - 2$ ,  $3 \cdot 2^{-2^{17}} < |\hat{a}/\hat{b}| < 3 \cdot 2^{2^{17}-1}$ , and  $\text{mbits}(\text{pc}) = \mu$ . Then  $q$  is normal,  $\hat{q}$  is  $(\mu + 1)$ -exact and

$$|\hat{q} - \hat{a}/\hat{b}| < 2^{\min(\text{expo}(\hat{q}), \text{expo}(\hat{a}/\hat{b}))-\mu}.$$

*Proof.* Let  $\alpha = 2^{-M}$ ,  $\beta = 2^{\text{expo}(\hat{a}/\hat{b})}$ , and  $\epsilon = 3/2^{16}$ . We define a sequence of approximations  $\xi_i$  of  $\hat{a}/\hat{b}$  by

$$\xi_i = \begin{cases} \hat{x}_0 & \text{if } i = 0 \\ \xi_{i-1}(2 - \hat{b}\xi_{i-1}) & \text{if } i > 0. \end{cases}$$

Since  $\hat{a}$  and  $\hat{b}$  are positive, so are the  $\xi_i$ , as well as every product computed by *FPU-MUL*. By Lemmas 4.5 and 4.7,  $|1 - \hat{b}\xi_i| < \epsilon^{2^i}$  for all  $i$ . Thus,  $\hat{b}\xi_i < 1 + \epsilon^{2^i}$  and  $2 - \hat{b}\xi_i < 1 + \epsilon^{2^i}$ . We also have

$$\hat{a}\xi_i = (\hat{a}/\hat{b})(\hat{b}\xi_i) < 2\beta(1 + \epsilon^{2^i})$$

and

$$|\hat{a}/\hat{b} - \hat{a}\xi_i| = (\hat{a}/\hat{b})|1 - \hat{b}\xi_i| < (\hat{a}/\hat{b})\epsilon^{2^i} < 2\beta\epsilon^{2^i}.$$

By Theorem 1,  $\hat{d}_0 = \text{near}(\hat{b}\hat{x}_0, M) = \text{near}(\hat{b}\xi_0, M)$ ; hence by Lemma 2.26,

$$|\hat{d}_0 - \hat{b}\xi_0| \leq 2^{\text{expo}(\hat{b}\xi_0)-M} \leq 2^{-M} = \alpha.$$

Note that our bounds for  $|\hat{a}/\hat{b}|$  ensure that the hypotheses of Theorem 1 are satisfied by  $x = a$  and  $y = x_0$ . Thus,

$$|\hat{n}_0 - \hat{a}\xi_0| \leq 2^{\text{expo}(\hat{a}\xi_0)-M} \leq 2^{\text{expo}(\hat{a}/\hat{b})+1-M} = 2\alpha\beta,$$

and by Lemma 3.5 (the hypotheses of which are ensured by Lemma 4.5),

$$0 < 2 - \hat{b}\xi_0 - 2\alpha \leq \hat{r}_0 < 2 - \hat{b}\xi_0.$$

Therefore,

$$\begin{aligned} \hat{n}_0\hat{r}_0 &< (\hat{a}\xi_0 + 2\alpha\beta)(2 - \hat{b}\xi_0) = \hat{a}\xi_1 + 2\alpha\beta(2 - \hat{b}\xi_0) < \hat{a}\xi_1 + 2\alpha\beta(1 + \epsilon) \\ &< \hat{a}\xi_1 + 2\alpha\beta + 2^{-13}\alpha\beta, \end{aligned}$$

$$\begin{aligned} \hat{n}_0\hat{r}_0 &\geq (\hat{a}\xi_0 - 2\alpha\beta)(2 - \hat{b}\xi_0 - 2\alpha) = \hat{a}\xi_1 - 2\alpha\beta(2 - \hat{b}\xi_0) - 2\alpha\hat{a}\xi_0 + 4\alpha^2\beta \\ &> \hat{a}\xi_1 - 2\alpha\beta(1 + \epsilon) - 2\alpha^2\beta(1 + \epsilon) > \hat{a}\xi_1 - 6\alpha\beta - 2^{-12}\alpha\beta, \end{aligned}$$

and

$$\begin{aligned} |\hat{n}_0\hat{r}_0 - \hat{a}/\hat{b}| &\leq |\hat{n}_0\hat{r}_0 - \hat{a}\xi_1| + |\hat{a}\xi_1 - \hat{a}/\hat{b}| < 7\alpha\beta + 2\beta\epsilon^2 \\ &< (7 \cdot 2^{-75} + 9 \cdot 2^{-31})\beta < 2^{-27}\beta \\ &= 2^{\text{expo}(\hat{a}/\hat{b})-27}. \end{aligned}$$

Suppose  $\text{pc} = \text{PC-32}$ . Then  $\mu = 24$  and

$$|\hat{n}_0\hat{r}_0 - \hat{a}/\hat{b}| < 2^{\text{expo}(\hat{a}/\hat{b})-27} < 2^{\text{expo}(\hat{a}/\hat{b})-\mu-2}.$$

By Lemma 3.7,  $\hat{q}$  is  $(\mu + 1)$ -exact,  $\text{expo}(\hat{n}_0\hat{r}_0) \leq \text{expo}(\hat{q})$ , and  $|\hat{n}_0\hat{r}_0 - \hat{q}| \leq 2^{\text{expo}(\hat{n}_0\hat{r}_0)-\mu-1}$ . We may now invoke Lemma 4.8 with  $\zeta = \hat{n}_0\hat{r}_0$ , which yields the desired inequality.

Thus, we may assume that  $\text{pc} \neq \text{PC-32}$ . Now

$$\hat{d}_0\hat{r}_0 < (\hat{b}\xi_0 + \alpha)(2 - \hat{b}\xi_0) = \hat{b}\xi_1 + \alpha(2 - \hat{b}\xi_0) < \hat{b}\xi_1 + \alpha + 2^{-14}\alpha,$$

$$\begin{aligned} \hat{d}_0\hat{r}_0 &\geq (\hat{b}\xi_0 - \alpha)(2 - \hat{b}\xi_0 - 2\alpha) = \hat{b}\xi_1 - 2\alpha\hat{b}\xi_0 - \alpha(2 - \hat{b}\xi_0) + 2\alpha^2 \\ &> \hat{b}\xi_1 - 2\alpha(1 + \epsilon) - \alpha(1 + \epsilon) > \hat{b}\xi_1 - 3\alpha - 2^{-13}\alpha, \end{aligned}$$

and Lemma 2.26 implies

$$|\hat{d}_1 - \hat{d}_0\hat{r}_0| \leq 2^{\text{expo}(\hat{d}_0\hat{r}_0)-M} \leq \alpha;$$

hence

$$\hat{d}_1 \leq \hat{d}_0\hat{r}_0 + \alpha < \hat{b}\xi_1 + 2\alpha + 2^{-14}\alpha$$

and

$$\hat{d}_1 \geq \hat{d}_0\hat{r}_0 - \alpha > \hat{b}\xi_1 - 4\alpha - 2^{-13}\alpha.$$

By Lemmas 3.5 and 3.6,

$$\hat{r}_1 < 2 - \hat{d}_0\hat{r}_0 < (2 - \hat{b}\xi_1) + 3\alpha + 2^{-13}\alpha$$

and

$$\hat{r}_1 \geq 2 - \hat{d}_0\hat{r}_0 - 2\alpha > (2 - \hat{b}\xi_1) - 3\alpha - 2^{-14}\alpha > 0.$$

Continuing in this manner, we have

$$|\hat{n}_1 - \hat{n}_0\hat{r}_0| \leq 2^{\text{expo}(\hat{n}_0\hat{r}_0)-M} \leq 2\alpha\beta,$$

$$\hat{n}_1 \leq \hat{n}_0\hat{r}_0 + 2\alpha\beta < \hat{a}\xi_1 + 4\alpha\beta + 2^{-13}\alpha\beta,$$

$$\hat{n}_1 \geq \hat{n}_0\hat{r}_0 - 2\alpha\beta > \hat{a}\xi_1 - 8\alpha\beta - 2^{-12}\alpha\beta,$$

$$\begin{aligned} \hat{n}_1\hat{r}_1 &< (\hat{a}\xi_1 + 4\alpha\beta + 2^{-13}\alpha\beta)((2 - \hat{b}\xi_1) + 3\alpha + 2^{-13}\alpha) \\ &< \hat{a}\xi_2 + (4\alpha\beta + 2^{-13}\alpha\beta)(1 + \epsilon^2) + 2\beta(1 + \epsilon^2)(3\alpha + 2^{-13}\alpha) \\ &\quad + (4\alpha\beta + 2^{-13}\alpha\beta)(3\alpha + 2^{-13}\alpha) \\ &< \hat{a}\xi_2 + 10\alpha\beta + 2^{-11}\alpha\beta, \end{aligned}$$

$$\begin{aligned} \hat{n}_1\hat{r}_1 &> (\hat{a}\xi_1 - 8\alpha\beta - 2^{-12}\alpha\beta)((2 - \hat{b}\xi_1) - 3\alpha + 2^{-14}\alpha) \\ &> \hat{a}\xi_2 - (8\alpha\beta + 2^{-12}\alpha\beta)(1 + \epsilon^2) - 2\beta(1 + \epsilon^2)(3\alpha + 2^{-14}\alpha) \\ &> \hat{a}\xi_2 - 14\alpha\beta - 2^{-11}\alpha\beta, \end{aligned}$$



and

$$\begin{aligned} |\hat{n}_1 \hat{r}_1 - \hat{a}/\hat{b}| &\leq |\hat{n}_1 \hat{r}_1 - \hat{a}\xi_2| + |\hat{a}\xi_2 - \hat{a}/\hat{b}| < 15\alpha\beta + 2^{-11}\alpha\beta + 2\beta\epsilon^4 \\ &< (15 \cdot 2^{-75} + 81 \cdot 2^{-63})\beta < 2^{-56}\beta \\ &= 2^{\text{expo}(\hat{a}/\hat{b})-56}. \end{aligned}$$

Suppose pc = PC-64, and therefore  $\mu = 53$ . Then

$$|\hat{n}_1 \hat{r}_1 - \hat{a}/\hat{b}| < 2^{\text{expo}(\hat{a}/\hat{b})-56} < 2^{\text{expo}(\hat{a}/\hat{b})-\mu-2}.$$

The remaining hypotheses of Lemma 4.8, with  $\hat{n}_1 \hat{r}_1$  substituted for  $\zeta$ , again follow from Lemma 3.7, and the desired inequality follows.

Thus, we may assume pc = PC-80 or pc = PC-87. Continuing, we have

$$\begin{aligned} \hat{d}_1 \hat{r}_1 &< (\hat{b}\xi_1 + 2\alpha + 2^{-14}\alpha)((2 - \hat{b}\xi_1) + 3\alpha + 2^{-13}\alpha) \\ &< \hat{b}\xi_2 + (2\alpha + 2^{-14}\alpha)(1 + \epsilon^2) + (3\alpha + 2^{-13}\alpha)(1 + \epsilon^2) \\ &\quad + (2\alpha + 2^{-14}\alpha)(3\alpha + 2^{-13}\alpha) \\ &< \hat{b}\xi_2 + 5\alpha + 2^{-12}\alpha, \end{aligned}$$

$$\begin{aligned} \hat{d}_1 \hat{r}_1 &> (\hat{b}\xi_1 - 4\alpha - 2^{-13}\alpha)((2 - \hat{b}\xi_1) - 3\alpha - 2^{-14}\alpha) \\ &> \hat{b}\xi_2 - (4\alpha + 2^{-13}\alpha)(1 + \epsilon^2) - (1 + \epsilon^2)(3\alpha + 2^{-14}\alpha) \\ &> \hat{b}\xi_2 - 7\alpha - 2^{-12}\alpha, \end{aligned}$$

$$\hat{r}_2 < 2 - \hat{d}_1 \hat{r}_1 < (2 - \hat{b}\xi_2) + 7\alpha + 2^{-12}\alpha,$$

$$\hat{r}_2 \geq 2 - \hat{d}_1 \hat{r}_1 - 2\alpha > (2 - \hat{b}\xi_2) - 7\alpha - 2^{-12}\alpha > 0,$$

$$|\hat{n}_2 - \hat{n}_1 \hat{r}_1| \leq 2^{\text{expo}(\hat{n}_1 \hat{r}_1)-M} \leq 2\alpha\beta,$$

$$\hat{n}_2 \leq \hat{n}_1 \hat{r}_1 + 2\alpha\beta < \hat{a}\xi_2 + 12\alpha\beta + 2^{-11}\alpha\beta,$$

$$\hat{n}_2 \geq \hat{n}_1 \hat{r}_1 - 2\alpha\beta > \hat{a}\xi_2 - 16\alpha\beta - 2^{-11}\alpha\beta,$$

$$\begin{aligned} \hat{n}_2 \hat{r}_2 &< (\hat{a}\xi_2 + 12\alpha\beta + 2^{-11}\alpha\beta)((2 - \hat{b}\xi_2) + 7\alpha + 2^{-12}\alpha) \\ &< \hat{a}\xi_3 + (12\alpha\beta + 2^{-11}\alpha\beta)(1 + \epsilon^4) + 2\beta(1 + \epsilon^4)(7\alpha + 2^{-12}\alpha) \\ &\quad + (12\alpha\beta + 2^{-11}\alpha\beta)(7\alpha + 2^{-12}\alpha) \\ &< \hat{a}\xi_3 + 26\alpha\beta + 2^{-9}\alpha\beta, \end{aligned}$$

and

$$\begin{aligned} \hat{n}_2 \hat{r}_2 &> (\hat{a}\xi_2 - 16\alpha\beta - 2^{-11}\alpha\beta)((2 - \hat{b}\xi_2) - 7\alpha + 2^{-12}\alpha) \\ &> \hat{a}\xi_3 - (16\alpha\beta + 2^{-11}\alpha\beta)(1 + \epsilon^4) - 2\beta(1 + \epsilon^4)(7\alpha + 2^{-12}\alpha) \\ &> \hat{a}\xi_3 - 30\alpha\beta - 2^{-9}\alpha\beta. \end{aligned}$$

Finally, since  $\mu \leq 68$ ,

$$\begin{aligned} |\hat{n}_2 \hat{r}_2 - \hat{a}/\hat{b}| &\leq |\hat{n}_2 \hat{r}_2 - \hat{a}\xi_3| + |\hat{a}\xi_3 - \hat{a}/\hat{b}| < 31\alpha\beta + 2\beta\epsilon^8 \\ &< (30 \cdot 2^{-75} + 81 \cdot 2^{-110})\beta < 2^{-70}\beta \\ &\leq 2^{\text{expo}(\hat{a}/\hat{b})-\mu-2}, \end{aligned}$$

and the lemma follows from Lemma 4.8, with  $\zeta = \hat{n}_2 \hat{r}_2$ . □

4.4. The operation OP-SQRT

The Newton-Raphson formula for approximating  $1/\sqrt{\hat{b}}$  is

$$\xi_i = \frac{\xi_{i-1}}{2} (3 - \hat{b} \xi_{i-1}^2).$$

Since the relative error of this approximation is

$$\left| \frac{\xi_i - 1/\sqrt{\hat{b}}}{1/\sqrt{\hat{b}}} \right| = |\sqrt{\hat{b}} \xi_i - 1| < |\sqrt{\hat{b}} \xi_i - 1| |\sqrt{\hat{b}} \xi_i + 1| = |\hat{b} \xi_i^2 - 1|,$$

convergence is established by the following lemma, which is proved in [8].

**Lemma 4.10.** *Let  $b, x \in \mathbb{Q}$  with  $0 \leq bx^2 \leq 4$  and let  $y = \frac{x}{2}(3 - bx^2)$ . Then*

$$0 \leq 1 - by^2 \leq (1 - bx^2)^2.$$

We shall use Lemma 4.10 to derive an error estimate for  $q$  in the OP-SQRT case.

**Lemma 4.11.** *For all  $i \in \mathbb{N}$ , let  $\xi_i$  be defined by*

$$\xi_i = \begin{cases} \hat{x}_0 & \text{if } i = 0 \\ \frac{\xi_{i-1}}{2} (3 - \hat{b} \xi_{i-1}^2) & \text{if } i > 0, \end{cases}$$

and let  $\epsilon = 3/2^{16}$ . Assume that  $\hat{q} > 0$  and  $\hat{q}$  is  $(\mu + 1)$ -exact, where  $\mu \geq 24$ .

Let  $\ell, h \in \mathbb{Q}$  such that  $0 \leq \ell \leq h$  and  $\ell^2 \leq \hat{b} \leq h^2$ . Let  $\zeta, \eta \in \mathbb{Q}^+$  and  $i \in \mathbb{Z}^+$  such that

$$\text{expo}(\zeta) \leq \text{expo}(\hat{q}),$$

$$|\hat{q} - \zeta| \leq 2^{\text{expo}(\zeta) - \mu - 1},$$

$$|\hat{b} \xi_i - \zeta| < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor} \eta,$$

and

$$2\eta + 8\epsilon^{2^i} \leq 2^{-\mu - 1}.$$

Then

$$h > q - 2^{\min(\text{expo}(\hat{q}), \text{expo}(h)) - \mu}$$

and

$$\ell < q + 2^{\min(\text{expo}(\hat{q}), \text{expo}(\ell)) - \mu}.$$

*Proof.* By Lemmas 4.6 and 4.10,  $0 \leq 1 - \hat{b} \xi_i^2 < \epsilon^{2^i}$ , where  $\epsilon = 3/2^{16}$ , and hence

$$(\hat{b} \xi_i)^2 = \hat{b} (\hat{b} \xi_i^2) > \hat{b} (1 - \epsilon^{2^i}) > 2^{\text{expo}(\hat{b}) - 1} > (2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - 1})^2$$

and  $\hat{b} \xi_i > 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - 1}$ .

Since  $|\hat{q} - \zeta| \leq 2^{\text{expo}(\zeta) - \mu - 1} \leq \zeta/4$ ,  $\hat{q} \geq \frac{3}{4}\zeta$ . Since  $\eta < 2^{-\mu - 2}$ ,

$$|\hat{b} \xi_i - \zeta| < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - \mu - 2} < \hat{b} \xi_i 2^{-\mu - 1} \leq \hat{b} \xi_i / 4,$$

and hence  $\hat{q} \geq \frac{3}{4}\zeta > \frac{9}{16}\hat{b}\xi_i$ , which implies

$$\hat{q}^2 > \frac{81}{256}(\hat{b}\xi_i)^2 > \frac{81}{256}\hat{b}(1 - \epsilon^{2^i}) > \hat{b}/4.$$

It follows that  $\text{expo}(\hat{q}) \geq \lfloor \text{expo}(\hat{b})/2 \rfloor - 1$ .

Since  $h^2 \geq \hat{b} \geq \hat{b}(\hat{b}\xi_i^2) = (\hat{b}\xi_i)^2$ ,

$$\begin{aligned} h &\geq \hat{b}\xi_i \geq \hat{q} - (|\hat{q} - \zeta| + |\hat{b}\xi_i - \zeta|) \\ &> \hat{q} - (2^{\text{expo}(\zeta) - \mu - 1} + 2^{\text{expo}(\hat{q}) - \mu - 1}) \geq \hat{q} - 2^{\text{expo}(\hat{q}) - \mu}. \end{aligned}$$

Therefore, we may assume  $\text{expo}(h) < \text{expo}(\hat{q})$ . But  $|h| > |\hat{q}|/2$ ; hence  $\text{expo}(h) = \text{expo}(\hat{q}) - 1$ . Also note that  $\text{expo}(h) \geq \lfloor \text{expo}(\hat{b})/2 \rfloor$ , for otherwise  $h < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor}$  and

$$\hat{b} \leq h^2 < 2^{2\lfloor \text{expo}(\hat{b})/2 \rfloor} \leq 2^{\text{expo}(\hat{b})}.$$

We may further assume  $\text{expo}(\zeta) = \text{expo}(\hat{q})$ , for otherwise  $\text{expo}(\zeta) \leq \text{expo}(h)$  and

$$\begin{aligned} h &\geq \hat{b}\xi_i \geq \hat{q} - (|\hat{q} - \zeta| + |\hat{b}\xi_i - \zeta|) \\ &> \hat{q} - (2^{\text{expo}(\zeta) - \mu - 1} + 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - \mu - 2}) \geq \hat{q} - 2^{\text{expo}(h) - \mu}. \end{aligned}$$

If  $\hat{q} > 2^{\text{expo}(\hat{q})}$ , then  $\hat{q} \geq 2^{\text{expo}(\hat{q})} + 2^{\text{expo}(\hat{q}) - \mu}$  by Lemma 2.13, and

$$h > \hat{q} - 2^{\text{expo}(\hat{q}) - \mu} \geq 2^{\text{expo}(\hat{q})} = 2^{\text{expo}(h) + 1}.$$

Therefore,  $\hat{q} = 2^{\text{expo}(\hat{q})}$ , which implies  $\zeta \geq \hat{q}$  and

$$h \geq \hat{b}\xi_i = \hat{q} - (\hat{q} - \hat{b}\xi_i) \geq \hat{q} - (\zeta - \hat{b}\xi_i) > \hat{q} - 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - \mu - 2} \geq \hat{q} - 2^{\text{expo}(h) - \mu}.$$

In order to derive the bound for  $\ell$ , we may assume  $\text{expo}(\hat{q}) \leq \text{expo}(\ell)$ , for otherwise  $\ell < \hat{q}$  and the inequality holds trivially. Since  $(\hat{b}\xi_i)^2 > \hat{b}(1 - \epsilon^{2^i})$ ,

$$\ell^2 \leq \hat{b} < (\hat{b}\xi_i)^2 / (1 - \epsilon^{2^i}) < [\hat{b}\xi_i / (1 - \epsilon^{2^i})]^2,$$

and hence

$$\ell < \hat{b}\xi_i / (1 - \epsilon^{2^i}) < \hat{b}\xi_i (1 + 2\epsilon^{2^i}).$$

Recall that  $\text{expo}(\hat{q}) \geq \lfloor \text{expo}(\hat{b})/2 \rfloor - 1$  and  $\hat{q} > \frac{9}{16}\hat{b}\xi_i$ ; hence  $\hat{b}\xi_i < 2^{\text{expo}(\hat{q}) + 2}$ . Thus,

$$\begin{aligned} \ell < \hat{b}\xi_i (1 + 2\epsilon^{2^i}) &< \hat{b}\xi_i + 8\epsilon^{2^i} 2^{\text{expo}(\hat{q})} \leq \hat{q} + |\hat{q} - \zeta| + |\zeta - \hat{b}\xi_i| + 8\epsilon^{2^i} 2^{\text{expo}(\hat{q})} \\ &< \hat{q} + 2^{\text{expo}(\hat{q}) - \mu - 1} + 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor} \eta + 8\epsilon^{2^i} 2^{\text{expo}(\hat{q})} \leq \hat{q} + 2^{\text{expo}(\hat{q})} (2^{-\mu - 1} + 2\eta + 8\epsilon^{2^i}) \\ &\leq \hat{q} + 2^{\text{expo}(\hat{q})} (2^{-\mu - 1} + 2^{-\mu - 1}) = \hat{q} + 2^{\text{expo}(\hat{q}) - \mu}. \end{aligned}$$

□

We shall also require the following lemma, in order to invoke Lemma 3.8.

**Lemma 4.12.** Under the hypothesis of Lemma 4.11,  $|\hat{q}^2 - \hat{b}| < 2^{\text{expo}(\hat{b}) - 3}$ .

*Proof.* Since  $\text{expo}(\hat{b}) \leq 2\lfloor \text{expo}(\hat{b})/2 \rfloor + 1$ ,  $\hat{b} < 2^{\text{expo}(\hat{b}) + 1} \leq (2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1})^2$ . Thus,

$$(\hat{b}\xi_i)^2 = \hat{b}(\hat{b}\xi_i^2) \leq \hat{b} < (2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1})^2$$

and  $\hat{b}\xi_i < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1}$ . Now since

$$|\hat{q} - \hat{b}\xi_i| \leq |\hat{q} - \zeta| + |\hat{b}\xi_i - \zeta| < 2^{\text{expo}(\hat{q}) - \mu} \leq 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1 - \mu}$$

and

$$|\hat{q} + \hat{b}\xi_i| \leq 2\hat{b}\xi_i + |\hat{q} - \hat{b}\xi_i| < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 2} + 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1 - \mu} < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 3},$$

we have

$$|\hat{q}^2 - (\hat{b}\xi_i)^2| = |\hat{q} - \hat{b}\xi_i| |\hat{q} + \hat{b}\xi_i| < 2^{2\lfloor \text{expo}(\hat{b})/2 \rfloor + 4 - \mu} \leq \hat{b}2^{4 - \mu}.$$

Thus,

$$|\hat{q}^2 - \hat{b}| \leq |\hat{q}^2 - (\hat{b}\xi_i)^2| + \hat{b}|1 - \hat{b}\xi_i^2| < \hat{b}2^{5 - \mu} < 2^{\text{expo}(\hat{b}) + 6 - \mu}.$$

□

**Lemma 4.13.** Assume  $\text{op} = \text{OP-SQRT}$ ,  $\hat{b} > 0$ ,  $\text{expo}(\hat{b}) \leq 2^{17} - 2$ , and let  $\text{mbits}(\text{pc}) = \mu$ . Let  $\ell, h \in \mathbb{Q}$  such that  $0 \leq \ell \leq h$  and  $\ell^2 \leq \hat{b} \leq h^2$ . Then  $q$  is normal,  $\hat{q}$  is  $(\mu + 1)$ -exact,

$$\ell < \hat{q} + 2^{\min(\text{expo}(\hat{q}), \text{expo}(\ell)) - \mu},$$

$$h > \hat{q} - 2^{\min(\text{expo}(\hat{q}), \text{expo}(h)) - \mu},$$

and

$$|\hat{q}^2 - \hat{b}| < 2^{\text{expo}(\hat{b}) - 3}.$$

*Proof.* Let  $\alpha = 2^{-M}$ ,  $\beta = 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor}$ , and  $\epsilon = 3/2^{16}$ . For  $i \in \mathbb{N}$ , let  $\xi_i$  be defined as in Lemma 4.11. Then  $\hat{b} < 4\beta^2$  and  $|1 - \hat{b}\xi_i^2| < \epsilon^{2^i}$ . For  $i > 0$ ,  $\hat{b}\xi_i^2 \leq 1$  and  $\hat{b}\xi_i < 2\beta$ , which implies  $2^{\text{expo}(\hat{b}\xi_i)} \leq \beta$ . On the other hand,

$$(\hat{b}\xi_0)^2 = \hat{b}(\hat{b}\xi_0^2) < 4\beta^2(1 + \epsilon) < (2\beta(1 + \epsilon))^2;$$

hence  $\hat{b}\xi_0 < 2\beta(1 + \epsilon) < 4\beta$ , which implies  $2^{\text{expo}(\hat{b}\xi_0)} \leq 2\beta$ . Also note that for all  $i$ ,

$$(3 - \hat{b}\xi_i^2)/2 = 1 + (1 - \hat{b}\xi_i^2)/2 < 1 + \epsilon^{2^i}/2.$$

We proceed as in the proof of Lemma 4.9, invoking Lemmas 4.11 and 4.12 in each of several cases. According to Lemma 4.6(c), the hypothesis of Theorem 1 is satisfied by  $x = y = \hat{x}_0$ . Thus,

$$\hat{t}_0 = \text{near}(\hat{x}_0^2, M) = \hat{x}_0^2 = \xi_0^2.$$

Similarly,

$$\hat{d}_0 = \text{near}(\hat{b}\xi_0, M)$$

and

$$\hat{n}_0 = \text{near}(\hat{b}\hat{t}_0, M) = \text{near}(\hat{b}\xi_0^2, M).$$

Therefore, by Lemma 2.26,

$$|\hat{d}_0 - \hat{b}\xi_0| \leq 2^{\text{expo}(\hat{b}\xi_0) - M} \leq 2\alpha\beta$$

and

$$|\hat{n}_0 - \hat{b}\xi_0^2| \leq 2^{\text{expo}(\hat{b}\xi_0^2) - M} \leq \alpha.$$

By Lemmas 3.5 and 4.6,

$$(3 - \hat{b}\xi_0^2)/2 - 2\alpha \leq \hat{r}_0 < (3 - \hat{b}\xi_0^2)/2.$$

Thus,

$$\hat{d}_0 \hat{r}_0 < (\hat{b}\xi_0 + 2\alpha\beta)(3 - \hat{b}\xi_0^2)/2 < \hat{b}\xi_1 + 2\alpha\beta(1 + \epsilon/2) < \hat{b}\xi_1 + 2\alpha\beta + 2^{-14}\alpha\beta$$

and

$$\begin{aligned} \hat{d}_0 \hat{r}_0 &> (\hat{b}\xi_0 - 2\alpha\beta)((3 - \hat{b}\xi_0^2)/2 - 2\alpha) > \hat{b}\xi_1 - 2\alpha\beta(1 + \epsilon/2) - 2\beta(1 + \epsilon)2\alpha \\ &> \hat{b}\xi_1 - 6\alpha\beta - 2^{-12}\alpha\beta. \end{aligned}$$

Suppose  $pc = PC-32$  and  $\mu = 24$ . We shall apply Lemmas 4.11 and 4.12 with  $\zeta = \hat{d}_0 \hat{r}_0$ ,  $i = 1$ , and  $\eta = 7\alpha$ . Under these substitutions, we have

$$|\hat{b}\xi_i - \zeta| < 7\alpha\beta = 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor} \eta$$

and

$$2\eta + 8\epsilon^{2^i} = 14 \cdot 2^{-M} + 8\epsilon^2 \leq 14 \cdot 2^{-75} + 9 \cdot 2^{-29} < 2^{-25} = 2^{-\mu-1}.$$

The remaining hypotheses of Lemma 4.11 are ensured by Lemma 3.7, and the conclusion follows.

Thus, we may assume  $pc \neq PC-32$ . Now we have  $\hat{t}_1 = \text{near}(\hat{r}_0^2, M)$ ; hence  $|\hat{t}_1 - \hat{r}_0^2| \leq \alpha$ , which implies

$$\hat{t}_1 \leq (3 - \hat{b}\xi_0^2)^2/4 + \alpha$$

and

$$\begin{aligned} \hat{t}_1 &\geq ((3 - \hat{b}\xi_0^2)/2 - 2\alpha)^2 - \alpha \\ &> (3 - \hat{b}\xi_0^2)^2/4 - 4\alpha(1 + \epsilon/2) - \alpha \\ &> (3 - \hat{b}\xi_0^2)^2/4 - 5\alpha - 2^{-13}\alpha. \end{aligned}$$

Consequently,

$$\begin{aligned} \hat{n}_0 \hat{t}_1 &\leq (\hat{b}\xi_0^2 + \alpha)((3 - \hat{b}\xi_0^2)^2/4 + \alpha) < \hat{b}\xi_1^2 + (1 + \epsilon)\alpha + \alpha(1 + \epsilon/2)^2 + \alpha^2 \\ &< \hat{b}\xi_1^2 + 2\alpha + 2^{-13}\alpha \end{aligned}$$

and

$$\begin{aligned} \hat{n}_0 \hat{t}_1 &\geq (\hat{b}\xi_0^2 - \alpha)((3 - \hat{b}\xi_0^2)^2/4 - 5\alpha - 2^{-13}\alpha) \\ &> \hat{b}\xi_1^2 - (1 + \epsilon)(5\alpha + 2^{-13}\alpha) - \alpha(1 + \epsilon/2)^2 > \hat{b}\xi_1^2 - 6\alpha - 2^{-11}\alpha. \end{aligned}$$

Since  $\hat{d}_1 = \text{near}(\hat{d}_0 \hat{r}_0, M)$ ,  $|\hat{d}_1 - \hat{d}_0 \hat{r}_0| \leq 2^{\text{expo}(\hat{d}_0 \hat{r}_0) - M} \leq 2\alpha\beta$ ; hence

$$\hat{b}\xi_1 - 8\alpha\beta - 2^{-12}\alpha\beta < \hat{d}_1 < \hat{b}\xi_1 + 4\alpha\beta + 2^{-14}\alpha\beta.$$

Similarly,  $\hat{n}_1 = \text{near}(\hat{n}_0 \hat{t}_1, M)$ ,  $|\hat{n}_1 - \hat{n}_0 \hat{t}_1| \leq 2^{\text{expo}(\hat{n}_0 \hat{t}_1) - M} \leq \alpha$ , and

$$\hat{b}\xi_1^2 - 7\alpha - 2^{-11}\alpha < \hat{n}_1 < \hat{b}\xi_1^2 + 3\alpha + 2^{-13}\alpha.$$

By Lemmas 3.5 and 3.6,

$$\hat{r}_1 < (3 - \hat{n}_0 \hat{t}_1)/2 < (3 - \hat{b}\xi_1^2)/2 + 3\alpha + 2^{-10}\alpha$$

and

$$\hat{r}_1 \geq (3 - \hat{n}_0 \hat{t}_1)/2 - 2\alpha > (3 - \hat{b}\xi_1^2)/2 - 3\alpha - 2^{-12}\alpha.$$

Thus,

$$\begin{aligned} \hat{d}_1 \hat{r}_1 &< (\hat{b}\xi_1 + 4\alpha\beta + 2^{-14}\alpha\beta)((3 - \hat{b}\xi_1^2)/2 + 3\alpha + 2^{-10}\alpha) \\ &< \hat{b}\xi_2 + 2\beta(3\alpha + 2^{-10}\alpha) + (4\alpha\beta + 2^{-14}\alpha\beta)(1 + \epsilon^2/2) \\ &\quad + (3\alpha + 2^{-10}\alpha)(4\alpha\beta + 2^{-14}\alpha\beta) \\ &< \hat{b}\xi_2 + 10\alpha\beta + 2^{-8}\alpha\beta \end{aligned}$$

and

$$\begin{aligned} \hat{d}_1 \hat{r}_1 &> (\hat{b}\xi_1 - 8\alpha\beta - 2^{-12}\alpha\beta)((3 - \hat{b}\xi_1^2)/2 - 3\alpha - 2^{-12}\alpha) \\ &> \hat{b}\xi_2 - 2\beta(3\alpha + 2^{-12}\alpha) - (8\alpha\beta + 2^{-12}\alpha\beta)(1 + \epsilon^2/2) \\ &> \hat{b}\xi_2 - 14\alpha\beta - 2^{-10}\alpha\beta. \end{aligned}$$

Suppose  $pc = PC-64$  and  $\mu = 53$ . We shall again invoke Lemmas 4.11 and 4.12, now with  $\zeta = \hat{d}_1 \hat{r}_1$ ,  $i = 2$ , and  $\eta = 15\alpha$ . Thus

$$|\hat{b}\xi_i - \zeta| < 15\alpha\beta = 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor} \eta$$

and

$$2\eta + 8\epsilon^{2^i} = 30 \cdot 2^{-M} + 8\epsilon^4 \leq 30 \cdot 2^{-75} + 81 \cdot 2^{-61} < 2^{-54} = 2^{-\mu-1}.$$

The remaining hypotheses of Lemma 4.11 are again ensured by Lemma 3.7.

Thus, we may assume  $pc = PC-80$  or  $pc = PC-87$ . Continuing in the same manner, we have

$$|\hat{t}_2 - \hat{r}_1^2| \leq 2^{\text{expo}(\hat{r}_1^2) - M} \leq \alpha,$$

$$\begin{aligned} \hat{t}_2 &< (3 - \hat{b}\xi_1^2)^2/4 + 2(1 + \epsilon^2/2)^2(3\alpha + 2^{-10}\alpha) + (3\alpha + 2^{-10}\alpha)^2 + \alpha \\ &< (3 - \hat{b}\xi_1^2)^2/4 + 7\alpha + 2^{-8}\alpha, \end{aligned}$$

$$\begin{aligned} \hat{t}_2 &> (3 - \hat{b}\xi_1^2)^2/4 - 2(1 + \epsilon^2/2)^2(3\alpha + 2^{-12}\alpha) - \alpha \\ &> (3 - \hat{b}\xi_1^2)^2/4 - 7\alpha + 2^{-10}\alpha, \end{aligned}$$

$$|\hat{d}_2 - \hat{d}_1 \hat{r}_1| \leq 2^{\text{expo}(\hat{d}_1 \hat{r}_1) - M} \leq 2\alpha\beta,$$

$$\hat{b}\xi_2 - 16\alpha\beta - 2^{-10}\alpha\beta < \hat{d}_2 < \hat{b}\xi_2 + 12\alpha\beta + 2^{-8}\alpha\beta,$$

$$\begin{aligned} \hat{n}_1 \hat{t}_2 &< (\hat{b}\xi_1^2 + 3\alpha + 2^{-13}\alpha)((3 - \hat{b}\xi_1^2)^2/4 + 7\alpha + 2^{-8}\alpha) \\ &< \hat{b}\xi_2^2 + (7\alpha + 2^{-8}\alpha) + (1 + \epsilon^2/2)^2(3\alpha + 2^{-13}\alpha) \\ &\quad + (3\alpha + 2^{-13}\alpha)(7\alpha + 2^{-8}\alpha) \\ &< \hat{b}\xi_1^2 + 10\alpha + 2^{-7}\alpha, \end{aligned}$$

$$\begin{aligned} \hat{n}_1 \hat{t}_2 &> (\hat{b}\xi_1^2 - 7\alpha - 2^{-11}\alpha)((3 - \hat{b}\xi_1^2)^2/4 - 7\alpha + 2^{-10}\alpha) \\ &> \hat{b}\xi_2^2 - (7\alpha + 2^{-10}\alpha) - (1 + \epsilon^2/2)^2(7\alpha + 2^{-11}\alpha) \\ &> \hat{b}\xi_1^2 - 14\alpha - 2^{-9}\alpha, \end{aligned}$$

$$\hat{r}_2 < (3 - \hat{n}_1 \hat{t}_2)/2 < (3 - \hat{b}\xi_2^2)/2 + 7\alpha + 2^{-10}\alpha,$$

$$\hat{r}_2 \geq (3 - \hat{n}_1 \hat{t}_2)/2 - 2\alpha > (3 - \hat{b}\xi_2^2)/2 - 7\alpha - 2^{-8}\alpha,$$

$$\begin{aligned} \hat{d}_2 \hat{r}_2 &< (\hat{b}\xi_2 + 12\alpha\beta + 2^{-8}\alpha\beta)((3 - \hat{b}\xi_2^2)/2 + 7\alpha + 2^{-10}\alpha) \\ &< \hat{b}\xi_3 + 2\beta(7\alpha + 2^{-10}\alpha) + (1 + \epsilon^4/2)(12\alpha\beta + 2^{-8}\alpha\beta) \\ &\quad + (12\alpha\beta + 2^{-8}\alpha\beta)(7\alpha + 2^{-10}\alpha) \\ &< \hat{b}\xi_3 + 26\alpha\beta + 2^{-7}\alpha\beta, \end{aligned}$$

and

$$\begin{aligned} \hat{d}_2 \hat{r}_2 &> (\hat{b}\xi_2 - 16\alpha\beta - 2^{-10}\alpha\beta)((3 - \hat{b}\xi_2^2)/2 - 7\alpha + 2^{-8}\alpha) \\ &> \hat{b}\xi_3 - 2\beta(7\alpha + 2^{-8}\alpha) - (1 + \epsilon^4/2)(16\alpha\beta + 2^{-10}\alpha\beta) \\ &> \hat{b}\xi_3 - 30\alpha\beta - 2^{-6}\alpha\beta. \end{aligned}$$

Finally, we apply Lemmas 4.11 and 4.12 with  $\zeta = \hat{d}_2 \hat{r}_2$ ,  $i = 3$ , and  $\eta = 31\alpha$ . Thus,

$$|\hat{b}\xi_i - \zeta| < 31\alpha\beta = 2^{\lceil \text{expo}(\hat{b})/2 \rceil} \eta,$$

and since  $\mu \leq 68$ ,

$$2\eta + 8\epsilon^{2^i} = 62 \cdot 2^{-M} + 8\epsilon^8 \leq 62 \cdot 2^{-75} + 2^{-112} < 2^{-69} \leq 2^{-\mu-1}.$$

The proof is completed by invoking Lemmas 3.7 and 4.11. □

#### 4.5. Final rounding

The remaining analysis pertains to the latter part of *FPU-DIV-SQRT*, in which the approximation  $q$  is adjusted to produce the correctly rounded result.

The significance of the variables  $q$ -guard and  $q$ -lsb is given by the following.

**Lemma 4.14.** Assume that  $q$  is normal and  $\hat{q}$  is  $(\mu + 1)$ -exact, where  $\mu = \text{mbits}(\text{pc})$ .

- (a)  $q$ -guard = 0  $\Leftrightarrow \hat{q}$  is  $\mu$ -exact;
- (b)  $q$ -lsb = 0  $\Leftrightarrow \text{trunc}(\hat{q}, \mu)$  is  $(\mu - 1)$ -exact.

*Proof.* (a) Let  $m = \text{get-man}(q)$ . Then  $m$  is  $(\mu + 1)$ -exact, i.e.,

$$m2^{\mu - \text{expo}(m)} = m2^{\mu+1-M} \in \mathbb{Z}$$

and

$$q\text{-guard} = m[M - \mu - 1] = \text{rem}(\lfloor m2^{\mu+1-M} \rfloor, 2) = \text{rem}(m2^{\mu+1-M}, 2).$$

But

$$m \text{ is } \mu\text{-exact} \Leftrightarrow m2^{\mu-M} \in \mathbb{Z} \Leftrightarrow m2^{\mu+1-M} \text{ is even} \Leftrightarrow q\text{-guard} = 0.$$

- (b)  $q$ -lsb =  $m[M - \mu] = \text{rem}(\lfloor m2^{\mu-M} \rfloor, 2)$  and  $\text{trunc}(m, \mu) = \lfloor m2^{\mu-M} \rfloor 2^{M-\mu}$ . Thus,  $\text{trunc}(m, \mu)$  is  $(\mu - 1)$ -exact  $\Leftrightarrow \lfloor m2^{\mu-M} \rfloor 2^{M-\mu} 2^{(\mu-1)-1-(M-1)} = \lfloor m2^{\mu-M} \rfloor / 2 \in \mathbb{Z}$   
 $\Leftrightarrow \lfloor m2^{\mu-M} \rfloor$  is even  
 $\Leftrightarrow q\text{-lsb} = 0.$

□

The correctness proof for division will be based on the following.

**Lemma 4.15.** Let  $\mu = \text{mbits}(\text{pc})$ . Suppose  $q$  is normal,  $\hat{q}$  is  $(\mu + 1)$ -exact,  $\text{sign} = 0$ , and  $2^{1-2^{17}} < \hat{q} < 2^{2^{17}}(2 - 2^{1-\mu})$ . Let  $x \in \mathbb{Q}$  such that

- (a)  $|x - \hat{q}| < 2^{\min(\text{expo}(\hat{q}), \text{expo}(x)) - \mu}$ ;
- (b) if  $\text{rem-neg} = 1$ , then  $\hat{q} > x$ ;
- (c) if  $\text{rem-pos} = 1$ , then  $\hat{q} < x$ ;
- (d) if  $\text{rem-zero} = 1$ , then  $\hat{q} = x$ .

Then  $z$  is normal and  $\text{rnd}(x, \text{rc}, \text{pc}) = \hat{z}$ .

*Proof.* Note that the hypothesis implies that  $\hat{z} > 0$  and  $x > 0$ .

*Case 1.*  $\text{rc} = \text{RC-NEG}$  or  $\text{rc} = \text{RC-CHOP}$

In this case,  $\text{rnd}(x, \text{rc}, \text{pc}) = \text{trunc}(x, \mu)$ .

*Subcase 1.1.*  $q\text{-guard} = 0$  and  $\text{rem-neg} = 1$

By Lemma 4.14,  $\hat{q}$  is  $\mu$ -exact. Also,  $x < \hat{q}$ . By Lemma 2.24,

$$\text{get-man}(q) \ \& \ (2^M - 2^{M-\mu}) = \text{trunc}(\text{get-man}(q), \mu) = \text{get-man}(q).$$

If  $\text{get-man}(q) = 2^{M-1}$ , then  $\hat{q} = 2^{\text{expo}(\hat{q})}$ , where by hypothesis,  $\text{expo}(\hat{q}) > 1 - 2^{17}$ . In this case,  $\hat{z} = (2 - 2^{1-\mu})2^{\text{expo}(\hat{q})-1}$  and  $\text{expo}(\hat{z}) = \text{expo}(\hat{q}) - 1$ . In all other cases,  $\hat{q} \geq 2^{\text{expo}(\hat{q})} + 2^{1+\text{expo}(\hat{q})-\mu}$ ,  $\hat{z} = \hat{q} - 2^{1+\text{expo}(\hat{q})-\mu}$ ,  $\hat{z} \geq 2^{\text{expo}(\hat{q})}$ , and  $\text{expo}(\hat{z}) = \text{expo}(\hat{q})$ . In any case,  $\hat{z} + 2^{1+\text{expo}(\hat{z})-\mu} = \hat{q}$ . Since  $\text{trunc}(x, \mu) \leq x < \hat{q}$ ,  $\text{trunc}(x, \mu) \leq \hat{z}$  by Lemma 2.13. Also,  $\text{trunc}(x, \mu) \geq \hat{z}$ , for otherwise we would have  $x < \hat{z}$ ,  $\text{expo}(x) \leq \text{expo}(\hat{z})$ , and

$$x > \hat{q} - 2^{\text{expo}(x)-\mu} > \hat{q} - 2^{1+\text{expo}(\hat{z})-\mu} = \hat{z}.$$

*Subcase 1.2.*  $q\text{-guard} = 1$

In this case,  $\hat{q}$  is not  $\mu$ -exact, and  $\hat{z} = \text{trunc}(\hat{q}, \mu)$ . By Lemma 2.27,  $\hat{z} = \hat{q} - 2^{\text{expo}(\hat{q})-\mu}$ . Therefore,

$$\text{trunc}(x, \mu) \leq x < \hat{q} + 2^{\text{expo}(\hat{q})-\mu} = \hat{z} + 2^{\text{expo}(\hat{q})+1-\mu} = \hat{z} + 2^{\text{expo}(\hat{z})+1-\mu},$$

and hence  $\text{trunc}(x, \mu) \leq \hat{z}$ . But since  $x > \hat{q} - 2^{\text{expo}(\hat{q})-\mu} = \hat{z}$ ,  $\text{trunc}(x, \mu) \geq \text{trunc}(\hat{z}, \mu) = \hat{z}$ .

*Subcase 1.3.*  $q\text{-guard} = \text{rem-neg} = 0$

$\hat{q}$  is  $\mu$ -exact,  $x \geq \hat{q}$ , and  $\hat{z} = \text{trunc}(\hat{q}, \mu) = \hat{q}$ .

In this case,

$$\text{trunc}(x, \mu) \leq x < \hat{q} + 2^{\text{expo}(\hat{q})-\mu} = \hat{z} + 2^{\text{expo}(\hat{z})-\mu} < \hat{z} + 2^{\text{expo}(\hat{z})+1-\mu},$$

which implies  $\text{trunc}(x, \mu) \leq \hat{z}$ . But  $x \geq \hat{q} = \hat{z}$  implies  $\text{trunc}(x, \mu) \geq \hat{z}$ .

*Case 2.*  $\text{rc} = \text{RC-POS}$

In this case,  $\text{rnd}(x, \text{rc}, \text{pc}) = \text{away}(x, \mu)$ .

*Subcase 2.1.*  $q\text{-guard} = 1$

Here,  $\hat{q}$  is  $(\mu + 1)$ -exact but not  $\mu$ -exact. By the same reasoning as used in Subcase 1.1, we may show that

$$\hat{z} = \text{trunc}(\hat{q}, \mu) + 2^{\text{expo}(\hat{q})+1-\mu}.$$

But then by Lemma 2.27,

$$\hat{z} = \hat{q} - 2^{\text{expo}(\hat{q})-\mu} + 2^{\text{expo}(\hat{q})+1-\mu} = \hat{q} + 2^{\text{expo}(\hat{q})-\mu} = \text{away}(\hat{q}, \mu).$$

Since  $x < \hat{q} + 2^{\text{expo}(\hat{q})-\mu} = \hat{z}$ ,  $\text{away}(x, \mu) \leq \text{away}(\hat{z}, \mu) = \hat{z}$ . But  $x > \hat{q} - 2^{\text{expo}(\hat{q})-\mu} = \text{trunc}(\hat{q}, \mu)$ ; hence  $\text{away}(x, \mu) \geq \text{trunc}(\hat{q}, \mu) + 2^{\text{expo}(\hat{q})+1-\mu} = \hat{z}$ .

*Subcase 2.2.*  $q\text{-guard} = 0$  and  $\text{rem-pos} = 1$ .



In this case,  $\hat{q}$  is  $\mu$ -exact,  $\hat{q} < x$ , and

$$\hat{z} = \text{trunc}(\hat{q}, \mu) + 2^{\text{expo}(\hat{q})+1-\mu} = \hat{q} + 2^{\text{expo}(\hat{q})+1-\mu}.$$

Since  $x < \hat{z}$ ,  $\text{away}(x, \mu) \leq \text{away}(\hat{z}, \mu) = \hat{z}$ . But  $\text{away}(x, \mu) \geq x > \hat{q}$ , so  $\text{away}(x, \mu) \geq \hat{q} + 2^{\text{expo}(\hat{q})+1-\mu} = \hat{z}$ .

*Subcase 2.3.*  $q$ -guard = rem-pos = 0

$\hat{q}$  is  $\mu$ -exact,  $x \leq \hat{q}$ , and  $\hat{z} = \text{trunc}(\hat{q}, \mu) = \hat{q}$ . Thus,

$$\text{away}(x, \mu) \leq \text{away}(\hat{q}, \mu) = \hat{q} = \hat{z}.$$

Since  $x > q - 2^{\text{expo}(x)-\mu}$ ,  $\text{away}(x, \mu) \geq \text{near}(x, \mu) \geq \hat{q}$  by Lemma 2.28.

*Case 3.* rc = RC-NEAR and  $q$ -guard = 0

Here,  $\hat{q}$  is  $\mu$ -exact,  $\text{rnd}(x, \text{rc}, \text{pc}) = \text{near}(x, \mu)$ , and  $\hat{z} = \text{trunc}(\hat{q}, \mu) = \hat{q}$ .

Since  $x < \hat{q} + 2^{\text{expo}(\hat{q})-\mu}$  implies  $\text{near}(x, \mu) \leq \hat{q} = \hat{z}$  by Lemma 2.28(b). But since  $x > \hat{q} - 2^{\text{expo}(x)-\mu}$ ,  $\text{near}(x, \mu) \geq \hat{q}$  by Lemma 2.28(c).

*Case 4.* rc = RC-NEAR and  $q$ -guard = 1

In this case,  $\hat{q}$  is  $(\mu + 1)$ -exact but not  $\mu$ -exact. Let  $a = q - 2^{\text{expo}(\hat{q})-\mu}$  and  $b = q + 2^{\text{expo}(\hat{q})-\mu}$ . By Lemma 2.27,  $a = \text{trunc}(\hat{q}, \mu)$  and  $b = \text{away}(\hat{q}, \mu)$ .

*Subcase 4.1.* rem-pos = 1

In this case,  $\hat{z} = b$  and  $\hat{q} < x$ . Since  $x < \hat{q} + 2^{\text{expo}(\hat{q})-\mu} = b$ ,

$$\text{near}(x, \mu) \leq \text{near}(b, \mu) = b = \hat{z}.$$

But  $x > q = b - 2^{\text{expo}(\hat{q})-\mu} \geq b - 2^{\text{expo}(x)-\mu}$ ; hence  $\text{near}(x, \mu) \geq b$ .

*Subcase 4.2.* rem-neg = 1

In this case,  $\hat{z} = \text{trunc}(\hat{q}, \mu) = a$  and  $x < \hat{q}$ ; hence  $\text{near}(x, \mu) \leq a = \hat{z}$  by Lemma 2.28, and  $x > q - 2^{\text{expo}(\hat{q})-\mu} = a$  implies  $\text{near}(x, \mu) \geq \text{near}(a, \mu) = a$ .

*Subcase 4.3.* rem-zero = 1

Here,  $x = \hat{q}$ ; hence  $\text{near}(x, \mu) = \text{near}(\hat{q}, \mu)$ . We shall show  $\text{near}(\hat{q}, \mu) = \hat{z}$ . Note that by Lemma 2.29,  $\text{near}(\hat{q}, \mu)$  is  $(\mu - 1)$ -exact.

If  $q$ -lsb = 1, then  $\hat{z} = b$  and  $a = \text{trunc}(\hat{q}, \mu)$  is not  $(\mu - 1)$ -exact by Lemma 4.14. Thus,  $\text{near}(\hat{q}, \mu) \neq a$ , which implies  $\text{near}(\hat{q}, \mu) = b = \hat{z}$ .

If  $q$ -lsb = 0, then  $\hat{z} = a$ ,  $a$  is  $(\mu - 1)$ -exact by Lemma 4.14. It follows that  $b$  is not  $(\mu - 1)$ -exact, and hence  $\text{near}(\hat{q}, \mu) = a$ . □

We may now state the correctness theorem for division. Note that the bound on  $\text{expo}(\hat{b})$  is required by Lemma 4.4 and is therefore unavoidable. The other constraint states that  $\text{expo}(\hat{a}/\hat{b})$  may not assume either of the limiting values  $1 - 2^{17}$  and  $2^{17}$ . This is acceptable since the hardware would never be expected to return a value with either of those exponents. In particular, IEEE compliance only involves exponents that are accommodated by the 80-bit (64, 15) format.

**Theorem 2.** Assume  $\text{op} = \text{OP-DIV}$ ,  $\text{rc}$  is a rounding control specifier,  $\text{pc}$  is an external precision control specifier, and  $a$  and  $b$  are normal encodings such that  $\text{expo}(\hat{b}) \leq 2^{17} - 2$  and  $2 - 2^{17} \leq \text{expo}(\hat{a}/\hat{b}) \leq 2^{17} - 1$ . Then  $z$  is a normal encoding and

$$\hat{z} = \text{rnd}(\hat{a}/\hat{b}, \text{rc}, \text{pc}).$$

*Proof.* By the same reasoning that was used in the proof of Theorem 1, we may assume that  $\hat{a} > 0$  and  $\hat{b} > 0$ . We need only show that the hypotheses of Lemma 4.15 are satisfied by  $x = \hat{a}/\hat{b}$ .

First note that our hypothesis regarding  $\text{expo}(\hat{a}/\hat{b})$  yields the bounds on  $|\hat{a}/\hat{b}|$  that are required by Lemma 4.9, which implies that  $\hat{q}$  is  $(\mu + 1)$ -exact and

$$|\hat{q} - \hat{a}/\hat{b}| < 2^{\min(\text{expo}(\hat{q}), \text{expo}(\hat{a}/\hat{b})) - \mu}.$$

This in turn implies the bounds on  $\hat{q}$  that are required by Lemma 4.15, as well as  $\hat{q} > 0$ , and hence  $\text{get-sign}(q) = \text{sign} = 0$ .

Next, we apply Lemma 3.8 with  $x = b$ ,  $y = q$ ,  $d = a$ , and  $z = \text{rem}$ , which implies that

$$|\hat{a}/\hat{b}| > |\hat{q}| \Leftrightarrow |\hat{b}\hat{q}| < |\hat{a}| \Leftrightarrow \text{get-man}(\text{rem})[M - 2] = 1 \Leftrightarrow \text{rem-pos} = 1$$

and

$$\hat{a}/\hat{b} = \hat{q} \Leftrightarrow \hat{b}\hat{q} = \hat{a} \Leftrightarrow \text{get-man}(\text{rem})[M - 2 : 0] = \text{inexact} = 0 \Leftrightarrow \text{rem-zero} = 1.$$

But since exactly one of  $\text{rem-pos}$ ,  $\text{rem-zero}$ , and  $\text{rem-neg}$  is nonzero, it follows that

$$|\hat{a}/\hat{b}| < |\hat{q}| \Leftrightarrow \text{rem-neg} = 1,$$

and all hypotheses of Lemma 4.15 are satisfied. □

In order to prove our correctness result for square root, a modification of Lemma 4.16 will be required.

**Lemma 4.16.** *Let  $\mu = \text{mbits}(\text{pc})$ . Suppose  $q$  is normal,  $\hat{q}$  is  $(\mu + 1)$ -exact,  $\text{sign} = 0$ , and  $2^{1-2^{17}} < \hat{q} < 2^{2^{17}}(2 - 2^{1-\mu})$ . Let  $\ell, h \in \mathbb{Q}$  such that*

- (a)  $\ell - 2^{\min(\text{expo}(\hat{q}), \text{expo}(\ell)) - \mu} < \hat{q} < h + 2^{\min(\text{expo}(\hat{q}), \text{expo}(h)) - \mu}$ ;
- (b) if  $\text{rem-neg} = 1$ , then  $\hat{q} > \ell$ ;
- (c) if  $\text{rem-pos} = 1$ , then  $\hat{q} < h$ ;
- (d) if  $\text{rem-zero} = 1$ , then  $\ell \leq \hat{q} \leq h$ .

Then  $z$  is normal and  $\text{rnd}(\ell, \text{rc}, \text{pc}) \leq \hat{z} \leq \text{rnd}(h, \text{rc}, \text{pc})$ .

*Proof.* We shall prove the first inequality; the proof of the second is similar.

*Case 1. rem-neg = 1*

Since  $\ell < \hat{q}$ , we may find  $x$  such that  $\ell < x < \hat{q}$  and  $x > \hat{q} - 2^{\min(\text{expo}(\hat{q}), \text{expo}(x)) - \mu}$ . Then  $\text{rnd}(\ell, \text{rc}, \text{pc}) \leq \text{rnd}(x, \text{rc}, \text{pc})$ , but by Lemma 4.15,  $\text{rnd}(x, \text{rc}, \text{pc}) = \hat{z}$ .

*Case 2. rem-pos = 1*

Choose  $x$  so that  $\hat{q} < x < q + 2^{\min(\text{expo}(\hat{q}), \text{expo}(x)) - \mu}$  and  $x > \ell$ . Then  $\text{rnd}(\ell, \text{rc}, \text{pc}) \leq \text{rnd}(x, \text{rc}, \text{pc})$ , but by Lemma 4.15,  $\text{rnd}(x, \text{rc}, \text{pc}) = \hat{z}$ .

*Case 3. rem-zero = 1*

Let  $x = \hat{q}$ . Then  $\ell \leq x$ ; hence  $\text{rnd}(\ell, \text{rc}, \text{pc}) \leq \text{rnd}(x, \text{rc}, \text{pc})$ , but by Lemma 4.15,  $\text{rnd}(x, \text{rc}, \text{pc}) = \hat{z}$ . □

**Theorem 3.** *Assume  $\text{op} = \text{OP-SQRT}$ ,  $\text{rc}$  is a rounding control specifier,  $\text{pc}$  is an external precision control specifier, and  $b$  is a normal encoding such that  $\text{expo}(\hat{b}) \leq 2^{17} - 2$ . Let  $\ell, h \in \mathbb{Q}$  such that  $0 \leq \ell \leq h$  and  $\ell^2 \leq \hat{b} \leq h^2$ . Then  $z$  is a normal encoding and*

$$\text{rnd}(\ell, \text{rc}, \text{pc}) \leq \hat{z} \leq \text{rnd}(h, \text{rc}, \text{pc}).$$

*Proof.* It suffices to show that the hypotheses of Lemmas 4.16 are satisfied. First, by Lemma 4.13,  $\hat{q}$  is  $(\mu + 1)$ -exact,

$$\ell < \hat{q} + 2^{\min(\text{expo}(\hat{q}), \text{expo}(\ell)) - \mu},$$

and

$$h > \hat{q} - 2^{\min(\text{expo}(\hat{q}), \text{expo}(h)) - \mu}.$$

Substituting  $2^{\lfloor \text{expo}(\hat{b})/2 \rfloor}$  for  $\ell$  in the same lemma, we have

$$\hat{q} > 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor} - 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - \mu} > 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - 1} > 0;$$

hence

$$\text{get-sign}(q) = 0 = \text{sign}.$$

Similarly, substituting  $2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1}$  for  $h$  yields

$$\hat{q} < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1} + 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 1 - \mu} < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 2}.$$

Thus,

$$2^{2-2^{17}} < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor - 1} \leq \hat{q} < 2^{\lfloor \text{expo}(\hat{b})/2 \rfloor + 2} < 2^{2^{17}}.$$

Finally, we apply Lemma 3.8 with  $x = y = q$ ,  $d = b$ , and  $z = \text{rem}$ , which yields the following.

- (1) if  $\text{rem-neg} = 1$ , then  $\hat{q}^2 > \hat{b} \geq \ell^2$ ; hence  $\hat{q} > \ell$ ;
- (2) if  $\text{rem-pos} = 1$ , then  $\hat{q}^2 < \hat{b} \leq h^2$ ; hence  $\hat{q} < h$ ;
- (3) if  $\text{rem-zero} = 1$ , then  $\hat{q}^2 = \hat{b}$ ; hence  $\ell \leq \hat{q} \leq h$ .

Thus, all hypotheses of Lemmas 4.16 are satisfied. □

## 5. Conclusion

As noted in the introduction, the practical value of formal verification has been illustrated in this exercise by the detection of two design flaws. Both of these were in the definition of the procedure *FPU-MUL*, but neither affected the results of multiplication. One was an error in the specification of the parameter  $r$  in the rare case in which  $\text{overflow} = 0$  and  $\text{round-carryout-no-overflow} = 1$ , which would inevitably have led to erroneous quotients and square roots for certain inputs. The other was in the calculation of  $z$  in the *OP-BACK* case, and might have led to improper rounding of square roots, although we were unable to exhibit a concrete example of this behavior. It was not surprising that neither problem was exposed by traditional testing methods. Once they had been identified, however, both were easily corrected before the design was committed to silicon.

Aside from the correction of errors, formal analysis may also provide insight that allows improvements in the efficiency of a design. For example, while the multiplier that was originally presented to us had a width of 76 bits, we were able to show, by representing it as an indefinite parameter  $M$ , that this width could effectively be reduced to 75 bits without sacrificing the accuracy of any of the operations that the multiplier supports.

Although the functionality of a physical device cannot be absolutely guaranteed by the properties of a mathematical model, a realistic model can provide a fairly high level of confidence. In this case, our analysis was based on a register-transfer model, far less abstract than the hardware models that are typically used in formal verification of floating point algorithms. It must be noted, however, that the evidentiality of our mechanical verification depends on the accuracy of several stages of manual translation. The original C encoding of the design was translated by hand into a special-purpose hardware description language,

from which a gate-level implementation was eventually constructed. Meanwhile, our verification began with the pseudocode representation of the C program on which the lemmas and theorems of this paper are based. After detailed proofs of all of these results were derived informally (and this paper was essentially written), the pseudocode was translated into ACL2 along with the lemma statements. Finally, formal proofs of these statements were generated mechanically by guiding the ACL2 prover through each step of the informal proofs.

Obviously, our confidence in the final product would be enhanced if we could eliminate or mechanize any of the steps in these translations. This has been a focus of our more recent work: we have implemented a mechanical translator from AMD's hardware description language directly to the logic of ACL2, thereby reducing the possibility of human error in the formalization of hardware designs. In a report that is yet to be released, we describe the use of this translator in the mechanical verification of the AMD-K7 floating point adder.

Of course, a successful formal verification project requires a significant investment. The cost to AMD of the results presented here was five months of the author's time, divided approximately equally between writing the informal proofs and checking them mechanically. Much of this time, however, was spent developing general methods and results, especially the theory of floating point arithmetic presented in Section 2, which could be reused in any floating point verification effort. We have already applied the same results to several problems, and it is our hope that others will find them useful in similar projects. Thus, the ACL2 formalization of this theory is included in [Appendix B](#)

#### Appendix A. *Input to the ACL2 prover*

This appendix is available to subscribers to the journal at:  
<http://www.lms.ac.uk/jcm/1/lms98001/appendix-a/>.

#### Appendix B. *An ACL2 library of floating point arithmetic*

This appendix is available to subscribers to the journal at:  
<http://www.lms.ac.uk/jcm/1/lms98001/appendix-b/>.

#### *References*

1. S.F. ANDERSON, J.G. EARLE, R.E. GOLDSCHMIDT and D.M. POWERS, 'The IBM System/360 Model 91 Floating Point Execution Unit', *IBM Journal of Research and Development*, 11 (January 1967) 34-53. 179
2. R.S. BOYER and J. MOORE, *A computational logic handbook* (Academic Press, Boston, MA, 1988). 148
3. R.E. BRYANT, 'Verification of arithmetic functions with binary moment diagrams', Technical Report CMU-CS-94-160, School of Computer Science, Carnegie-Mellon University, 1994. 148
4. E.M. CLARKE and X. ZHAO, 'Word level symbolic model checking: a new approach for verifying arithmetic circuits', Technical Report CMU-CS-95-161, School of Computer Science, Carnegie-Mellon University, 1995. 148
5. INSTITUTE OF ELECTRICAL AND ELECTRONIC ENGINEERS, 'IEEE Standard for Binary Floating Point Arithmetic', Std. 754-1985, (IEEE, New York, NY, 1985). 148, 149

6. J. MOORE, T. LYNCH and M. KAUFMANN, 'A mechanically checked proof of the correctness of the kernel of the *AMD5<sub>K</sub>86* floating point division algorithm', *IEEE Transactions on Computers*, 47 (September, 1998). 148
7. S.F. OBERMAN, 'Division and square root for the AMD-K7 FPU' (Advanced Micro Devices, Milpitas, CA, March 1997). 149
8. D.M. RUSSINOFF, 'A mechanically checked proof of IEEE compliance of the AMD-K5 floating point square root microcode', *Formal Methods in System Design*, to appear. <http://www.onr.com/user/russ/david/fsqrt.html>. 148, 149, 150, 150, 157, 159, 189
9. G.L. STEELE, Jr., *Common Lisp The Language* 2nd edition (Digital Press, Waltham MA, 1990). 149

David M. Russinoff [david.russinoff@amd.com](mailto:david.russinoff@amd.com)

Advanced Micro Devices, Inc.  
5900 E. Ben White Blvd  
MS 625  
Austin, TX 78741  
U.S.A.