

## ACCURATE FLOATING-POINT SUMMATION \*

SIEGFRIED M. RUMP <sup>†</sup>, TAKESHI OGITA <sup>‡</sup>, AND SHIN'ICHI OISHI <sup>§</sup>

**Abstract.** Given a vector of floating-point numbers with exact sum  $s$ , we present an algorithm for calculating a faithful rounding of  $s$  into the set of floating-point numbers, i.e. one of the immediate floating-point neighbors of  $s$ . If the  $s$  is a floating-point number, we prove that this is the result of our algorithm. The algorithm adapts to the condition number of the sum, i.e. it is very fast for mildly conditioned sums with slowly increasing computing time proportional to the condition number. All statements are also true in the presence of underflow. Furthermore algorithms with  $K$ -fold accuracy are derived, where in that case the result is stored in a vector of  $K$  floating-point numbers. We also present an algorithm for rounding the sum  $s$  to the nearest floating-point number. Our algorithms are fast in terms of measured computing time because they neither require special operations such as access to mantissa or exponent, they contain no branch in the inner loop, nor do they require extra precision: The only operations used are standard floating-point addition, subtraction and multiplication in one working precision, for example double precision. Moreover, in contrast to other approaches, the algorithms are ideally suited for parallelization. We also sketch dot product algorithms with similar properties.

**Key words.** maximally accurate summation, faithful rounding, maximally accurate dot product, error-free transformations, fast algorithms, parallel algorithms, high precision

**AMS subject classifications.** 15-04, 65G99, 65-04

**1. Introduction.** We will present fast algorithms to compute approximations of high quality of the sum and the dot product of vectors of floating-point numbers. In the recent paper [36] we developed such algorithms delivering a result *as if* computed in a specified *precision*. Those algorithms use only one working precision, and we could generate a result *as if* computed in  $K$ -fold precision,  $K \geq 1$ . The *accuracy* of the computed result, however, still depends on the condition number of the problem.

In the present paper we go one step further and present algorithms producing a result of specified *accuracy*. Again using only floating-point operations in one working precision we will show that we can compute a floating-point result **res** such that there is no other floating-point number between the true, real result and **res**.

This has been called *faithful rounding* in the literature [10, 39, 9]. We will also prove that if the true result is an exactly representable floating-number, then this will be our computed result. This implies that the sign of the result is always determined correctly, a significant problem in the computation of geometrical predicates [19, 43, 7, 27, 6, 12], where the sign of the value of a dot product decides whether a point is exactly on a plane or on which side it is.

As we will show, the computational effort of our methods is proportional to the condition number of the problem. An almost ideal situation: for simple problems the algorithm is fast, and slows down with increasing difficulty.

---

\*This research was partially supported by Grant-in-Aid for Specially Promoted Research (No. 17002012: Establishment of Verified Numerical Computation) from the Ministry of Education, Science, Sports and Culture of Japan.

<sup>†</sup>Institute for Reliable Computing, Hamburg University of Technology, Schwarzenbergstraße 95, Hamburg 21071, Germany, and Visiting Professor at Waseda University, Faculty of Science and Engineering, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan ([rump@tu-harburg.de](mailto:rump@tu-harburg.de)).

<sup>‡</sup>CREST, Japan Science and Technology Agency (JST), and Faculty of Science and Engineering, Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan ([ogita@waseda.jp](mailto:ogita@waseda.jp)).

<sup>§</sup>Department of Computer Science, Faculty of Science and Engineering, Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan ([oishi@waseda.jp](mailto:oishi@waseda.jp)).

In addition to the faithfully rounded result in working precision, we also present algorithms for a rounded to nearest result. As is known, here the computational effort cannot be proportional to the traditional condition number of the problem, but depends on the nearness of the true result to a switching point for rounding.

Moreover, algorithms for producing a faithfully rounded result in  $K$ -fold accuracy are presented. In this case the result is represented as a vector of  $K$  floating-point numbers.

Our algorithms are fast. We interpret fast not only by the number of floating-point operations, but in terms of *measured computing time*. This means that special operations such as rounding to integer, access to mantissa or exponent, branches etc. are avoided. Our algorithms use only floating-point addition, subtraction and multiplication in working precision. No extra precision is required.

Summation and dot product are most basic tasks in numerical analysis, and there are numerous algorithms for that, among them [4, 8, 11, 12, 17, 21, 22, 23, 24, 25, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 44, 45]. Higham [18] devotes an entire chapter to summation. Accurate summation or dot product algorithms have various applications in many different areas of numerical analysis. Excellent overviews can be found in [18, 31].

Many known algorithms produce a 'more accurate' result, demonstrated by examples or rigorous estimations. There are few methods producing a result as if computed in higher precision and very few methods computing a result with specified accuracy [4, 39, 11, 28]. In [28], for example, the fact is used that IEEE 754 double precision numbers can be added exactly in some superlong accumulator due to their limited exponent range. With some bit manipulation the rounded to nearest result of the sum (or dot product) can be obtained.

For more details see [36]. However, known methods require sorting of the input vector, access to mantissa/exponent, branches, some extended precision and/or other special operations. As will be shown in the computational results, this may slow down a computation substantially.

An important tool in our work are so-called error-free transformations. Here a floating-point approximation is supplemented by an exact error term. Consider the following algorithm for the product of two floating-point numbers.

ALGORITHM 1.1. *Error-free transformation of the product of two floating-point numbers.*

```
function [x, y] = TwoProduct(a, b)
    x = fl(a · b)
    [a1, a2] = Split(a)
    [b1, b2] = Split(b)
    y = fl(a2 · b2 - (((x - a1 · b1) - a2 · b1) - a1 · b2))
```

Here  $\text{fl}(\cdot)$  denotes that the expression inside the parenthesis is calculated in floating-point. If no underflow occurs, this algorithm due to G.W. Veltkamp (see [10]) satisfies for all floating-point numbers  $a, b$

$$(1.1) \quad x = \text{fl}(a + b) \quad \text{and} \quad x + y = a \cdot b.$$

It relies on the ingenious splitting by Dekker [10] of a 53-bit floating-point number into two 26-bit parts. This is possible using the sign bit and allows exact multiplication of the parts.

ALGORITHM 1.2. *Error-free splitting of a floating-point number into two parts.*

```
function [x, y] = Split(a)
    c = fl(factor · a)    % factor = 227 + 1
    x = fl(c - (c - a))
    y = fl(a - x)
```

Such transformations were extensively used in our algorithms for summation and dot product [36] giving a result as if computed in a specified precision. Recently there is increasing interest in using such error-free transformations [14, 29].

By Algorithm 1.1 (`TwoProduct`) we can transform the dot product of two vectors of floating-point numbers into a sum of floating-point numbers, thus we concentrate on summation.

In the following we need many careful floating-point estimations, frequently heavily relying on bit representations and the definition of the floating-point arithmetic in use. Not only that this is frequently quite tedious, it is also sometimes presented in a colloquial manner and not easy to follow. To avoid this and also to ensure rigor, we found it convenient and more stringent to use inequalities. For this we derived some machinery to characterize floating-point numbers, their bit representations and to handle delicate situations.

The paper is organized as follows. First we introduce our notation in Section 2 and list a number of properties. Some of them may be trivial, however, we found them useful to note to get acquainted with the machinery. In this section we also define faithful rounding and give a sufficient criterion for it. In Section 3 we use this to develop an error-free transformation of a vector of floating-point numbers into an approximation of the sum and some remaining part. The magnitude of the remaining part can be estimated, so that we can introduce summation algorithms with faithful rounding in the following Section 4. We prove faithfulness which especially includes the exact determination of the sign. This is not only also true in the presence of underflow, but the computed result is exact if it is in the underflow range. We also estimate the computing time depending on the condition number. In the next Section 5 we extend the results to  $K$ -fold faithful rounding, where the result is represented as a vector of  $K$  floating-point numbers. In Section 6 we derive, based on the previous ones, a rounding to nearest algorithm. Up to now the vector lengths were restricted to about  $\sqrt{\text{eps}^{-1}}$ ; in the following Section 7 we extend the previous results to huge vector lengths up to almost  $\text{eps}^{-1}$ . A simplified and efficient version only for sign determination is given as well. Next we sketch algorithms for dot products with similar properties, and finally we give computational results in Section 9. In the Appendix we give some proof and executable Matlab code, and a Summary concludes the paper.

As in [36], all theorems, error analysis and proofs are due to the first author.

**2. Basic facts.** In this section we collect some basic facts for the analysis of our algorithms. Throughout the paper we assume that no overflow occurs, but we allow underflow. We will use only one working precision for all floating-point computations; as an example we sometimes refer to IEEE 754 double precision. This corresponds to 53 bits precision including an implicit bit for normalized numbers. However, we stress that the following analysis applies *mutatis mutandis* to a binary arithmetic with another format or exponent range by replacing the roundoff and underflow unit, for example to IEEE 754 single precision. Since we use floating-point numbers in only one working precision, we can refer to them as “the floating-point numbers”.

The set of all floating-point numbers is denoted by  $\mathbb{F}$ , and  $\mathbb{U}$  denotes the set of unnormalized floating-point numbers together with zero and the two normalized floating-point numbers of smallest nonzero magnitude. The relative rounding error unit, the distance from 1.0 to the next smaller<sup>1</sup> floating-point number, is denoted by  $\text{eps}$ , and the underflow unit by  $\text{eta}$ , that is the smallest positive (unnormalized) floating-point number. For IEEE 754 double precision we have  $\text{eps} = 2^{-53}$  and  $\text{eta} = 2^{-1074}$ . Then  $\frac{1}{2}\text{eps}^{-1}\text{eta}$  is the smallest positive normalized floating-point number and for  $f \in \mathbb{F}$  we have

$$(2.1) \quad f \in \mathbb{U} \Leftrightarrow 0 \leq 2\text{eps}|f| \leq \text{eta} .$$

<sup>1</sup>Note that sometimes the distance from 1.0 to the next *larger* floating-point number is used; for example, Matlab adopts this rule.

We denote by  $\text{fl}(\cdot)$  the result of a floating-point computation, where all operations within the parentheses are executed in working precision. If the order of execution is ambiguous and is crucial, we make it unique by using parentheses. We assume the floating-point operations to satisfy the properties of IEEE 754 arithmetic standard in rounding to nearest [20]. Then floating-point addition and subtraction satisfy [18]

$$(2.2) \quad \text{fl}(a \circ b) = (a \circ b)(1 + \varepsilon) \quad \text{for } a, b \in \mathbb{F}, \circ \in \{+, -\} \quad \text{and } |\varepsilon| \leq \mathbf{eps}.$$

Note that addition and subtraction is exact in case of underflow [15], so we need no underflow unit in (2.2). More precisely,

$$(2.3) \quad a, b \in \mathbb{F} \quad \text{and} \quad |a + b| \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta} \quad \text{implies} \quad \text{fl}(a + b) = a + b.$$

We have to distinguish between normalized and unnormalized floating-point numbers since the latter lack the implicit bit. As has been noted by several authors [34, 26, 10], the error of a floating-point addition is always a floating-point number:

$$(2.4) \quad a, b \in \mathbb{F} \quad \text{implies} \quad \delta := \text{fl}(a + b) - (a + b) \in \mathbb{F}.$$

Fortunately, the error term  $\delta$  can be computed using only standard floating-point operations. The following algorithm by Knuth was already given in 1969 [26].

**ALGORITHM 2.1.** *Error-free transformation for the sum of two floating-point numbers.*

```
function [x, y] = TwoSum(a, b)
    x = fl(a + b)
    z = fl(x - a)
    y = fl((a - (x - z)) + (b - z))
```

Knuth's algorithm transforms any pair of floating-point numbers  $(a, b)$  into a new pair  $(x, y)$  with

$$(2.5) \quad x = \text{fl}(a + b) \quad \text{and} \quad x + y = a + b.$$

This is also true in the presence of underflow. An error-free transformation for subtraction follows since  $\mathbb{F} = -\mathbb{F}$ . The  $\text{fl}(\cdot)$  notation applies not only to operations but to real numbers as well. For  $r \in \mathbb{R}$ ,  $\text{fl}(r) \in \mathbb{F}$  is  $r$  rounded to the nearest floating-point number. Following the IEEE 754 arithmetic standard tie is rounded to even. For  $f_1, f_2 \in \mathbb{F}$  and  $r \in \mathbb{R}$ , monotonicity of the rounding implies

$$(2.6) \quad f_1 \leq r \leq f_2 \quad \Rightarrow \quad f_1 \leq \text{fl}(r) \leq f_2$$

$$(2.7) \quad f_1 < \text{fl}(r) < f_2 \quad \Rightarrow \quad f_1 < r < f_2.$$

In numerical analysis the accuracy of a result is sometimes measured by the ‘‘unit in the last place (ulp)’’. For the following, sometimes delicate error estimations the ulp-concept has the drawback that it depends on the floating-point format and needs extra care in the underflow range.

We found it useful to introduce the ‘‘unit in the first place’’ (ufp) or leading bit of a real number by

$$(2.8) \quad 0 \neq r \in \mathbb{R} \quad \Rightarrow \quad \text{ufp}(r) := 2^{\lfloor \log_2 |r| \rfloor},$$

where we set  $\text{ufp}(0) := 0$ . This gives a convenient way to characterize the bits of a normalized floating-point number  $f$ : they range between the leading bit  $\text{ufp}(f)$  and the unit in the last place  $2\mathbf{eps} \cdot \text{ufp}(f)$ . In our analysis we will view a floating-number frequently as a scaled integer. For  $\sigma = 2^k, k \in \mathbb{Z}$ , we will use the set  $\mathbf{eps}\sigma\mathbb{Z}$ , which can be viewed as a set of fixed point numbers with smallest positive number  $\mathbf{eps}\sigma$ . The situation is depicted in Figure 2.1. Of course,  $\mathbb{F} \subseteq \mathbf{eta}\mathbb{Z}$ . Note that (2.8) is independent of some floating-

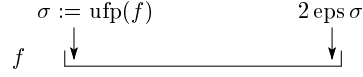


FIG. 2.1. Normalized floating-point number: unit in the first place and unit in the last place

point format and it applies to real numbers as well:  $\text{ufp}(r)$  is the value of the first nonzero bit in the binary representation of  $r$ . It follows

$$(2.9) \quad 0 \neq r \in \mathbb{R} \quad \Rightarrow \quad \text{ufp}(r) \leq |r| < 2\text{ufp}(r)$$

$$(2.10) \quad r, r' \in \mathbb{R} \text{ and } \text{ufp}(r) \leq |r'| \quad \Rightarrow \quad \text{ufp}(r) \leq \text{ufp}(r') .$$

We collect some properties. For  $\sigma = 2^k, k \in \mathbb{Z}, r \in \mathbb{R}$  it follows

$$(2.11) \quad \sigma' = 2^m, m \in \mathbb{Z} \text{ and } \sigma' \geq \sigma \quad \Rightarrow \quad \text{eps}\sigma'\mathbb{Z} \subseteq \text{eps}\sigma\mathbb{Z}$$

$$(2.12) \quad f \in \mathbb{F} \text{ and } |f| \geq \sigma \quad \Rightarrow \quad \text{ufp}(f) \geq \sigma$$

$$(2.13) \quad f \in \mathbb{F} \quad \Rightarrow \quad f \in 2\text{eps} \cdot \text{ufp}(f)\mathbb{Z}$$

$$(2.14) \quad r \in \text{eps}\sigma\mathbb{Z}, |r| \leq \sigma \text{ and } \text{eps}\sigma \geq \text{eta} \quad \Rightarrow \quad r \in \mathbb{F}$$

$$(2.15) \quad a, b \in \mathbb{F} \cap \text{eps}\sigma\mathbb{Z} \text{ and } \delta := \text{fl}(a+b) - (a+b) \quad \Rightarrow \quad \text{fl}(a+b), a+b, \delta \in \text{eps}\sigma\mathbb{Z}$$

$$(2.16) \quad a_i \in \mathbb{F} \text{ and } |a_i| \leq \sigma \quad \Rightarrow \quad \left| \text{fl}\left(\sum_{i=1}^n a_i\right) \right| \leq n\sigma$$

$$(2.17) \quad a, b \in \mathbb{F} \quad \Rightarrow \quad \begin{array}{l} \text{fl}(a+b) \in \text{eps} \cdot \text{ufp}(a)\mathbb{Z} \text{ and} \\ \text{fl}(a+b) \in \text{eps} \cdot \text{ufp}(b)\mathbb{Z} . \end{array}$$

The assertions (2.11) to (2.15) are clear. Note that (2.13) is also true for  $f \in \mathbb{U}$ . Estimation (2.16) is sometimes useful to avoid unnecessary quadratic terms and it is valid for any order of summation. It follows by induction: For  $s := \text{fl}(\sum_{i \neq k} a_i)$  we have  $|s + a_k| \leq n\sigma$ , so by (2.6) and if  $n\sigma \in \mathbb{F}$  also  $|\text{fl}(s + a_k)| \leq n\sigma$ , and if  $n\sigma$  is in the overflow range but  $\text{fl}(\sum a_i)$  is not, the estimation is true as well.

The last one (2.17) is also clear after a little thinking, and a rigorous proof follows easily with our machinery. The assertion is clear for  $ab \geq 0$  by using (2.13) and (2.11) because then  $|\text{fl}(a+b)| \geq \max(|a|, |b|)$ ; so without loss of generality it suffices to show  $\text{fl}(a-b) \in \text{eps}\sigma\mathbb{Z}$  for  $a \geq b \geq 0$  and  $\sigma := \text{ufp}(a)$ . If  $\text{ufp}(b) \geq \frac{1}{2}\sigma$ , then (2.13) implies  $a, b \in \text{eps}\sigma\mathbb{Z}$  and the assertion follows by (2.15). And if  $\text{ufp}(b) < \frac{1}{2}\sigma$ , then  $b < \frac{1}{2}\sigma$ , and  $a \geq \sigma$  implies  $a-b > \frac{1}{2}\sigma \in \mathbb{F}$ . Hence (2.6) shows  $\text{fl}(a-b) \geq \frac{1}{2}\sigma \in \mathbb{F}$  and (2.13) implies  $\text{fl}(a-b) \in \text{eps}\sigma\mathbb{Z}$ .

For later use we collect some more properties. For  $r \in \mathbb{R}$  and  $\tilde{r} = \text{fl}(r)$ ,

$$(2.18) \quad \tilde{r} \neq 0 \quad \Rightarrow \quad \text{ufp}(r) \leq \text{ufp}(\tilde{r})$$

$$(2.19) \quad \begin{array}{l} \tilde{r} \in \mathbb{F} \setminus \mathbb{U} \quad \Rightarrow \quad |\tilde{r} - r| \leq \text{eps} \cdot \text{ufp}(r) \leq \text{eps} \cdot \text{ufp}(\tilde{r}) \\ \tilde{r} \in \mathbb{U} \quad \Rightarrow \quad |\tilde{r} - r| \leq \frac{1}{2}\text{eta} . \end{array}$$

Note that a strict inequality occurs in (2.18) iff  $\tilde{r}$  is a power of 2 and  $|r| < |\tilde{r}|$ . The assertions follow by the rounding to nearest property of  $\text{fl}(\cdot)$ . Applying (2.19), (2.18), (2.9) and (2.3) to floating-point addition yields for  $a, b \in \mathbb{F}$ ,

$$(2.20) \quad f = \text{fl}(a+b) \quad \Rightarrow \quad f = a+b+\delta \text{ with } |\delta| \leq \text{eps} \cdot \text{ufp}(a+b) \leq \text{eps} \cdot \text{ufp}(f) \leq \text{eps}|f| .$$

We will frequently need this refined version of the standard estimation (2.2). Note that (2.20) is also true in the presence of underflow because in this case  $\delta = 0$ , the addition is exact. The  $\text{ufp}$  concept also allows

simple sufficient conditions for the fact that a floating-point addition is exact. For  $a, b \in \mathbb{F}$  and  $\sigma = 2^k$ ,  $k \in \mathbb{Z}$ ,

$$(2.21) \quad \begin{aligned} a, b \in \mathbf{eps}\sigma\mathbb{Z} \quad \text{and} \quad \text{fl}(|a+b|) < \sigma &\Rightarrow \text{fl}(a+b) = a+b \quad \text{and} \\ a, b \in \mathbf{eps}\sigma\mathbb{Z} \quad \text{and} \quad |a+b| \leq \sigma &\Rightarrow \text{fl}(a+b) = a+b . \end{aligned}$$

We need only to prove the second part since  $\text{fl}(|a+b|) < \sigma$  and (2.7) imply  $|a+b| < \sigma$ . To see the second part we first note that  $a+b \in \mathbf{eps}\sigma\mathbb{Z}$ . By (2.3) the addition is exact if  $|a+b| \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , and also if  $|a+b| = \sigma$ . Otherwise, (2.9) and (2.12) yield  $\sigma > |a+b| \geq \text{ufp}(a+b) \geq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  since  $\frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  is a power of 2, so  $\mathbf{eps}\sigma \geq 2\mathbf{eps} \cdot \text{ufp}(a+b) \geq \mathbf{eta}$  and (2.14) do the job.

The well-known result by Sterbenz [18, Theorem 2.5] says that subtraction is exact if floating-point numbers  $a, b \in \mathbb{F}$  of the same sign are not too far apart. More precisely, for  $a, b \geq 0$  we have

$$(2.22) \quad \frac{1}{2}a \leq b \leq 2a \quad \Rightarrow \quad \text{fl}(b-a) = b-a.$$

As an example of a proof we mention that with our machinery this is not difficult to see. If  $b \geq a$ , then (2.13) implies  $a, b, a-b \in 2\mathbf{eps}\sigma\mathbb{Z}$  for  $\sigma := \text{ufp}(a)$ . By assumption and (2.9),  $|b-a| = b-a \leq a < 2\sigma$ , and (2.21) proves this part. For  $b < a$ , (2.13) implies  $a, b, a-b \in 2\mathbf{eps}\sigma\mathbb{Z}$  for  $\sigma := \text{ufp}(b)$ , and similarly  $|b-a| = a-b \leq b < 2\sigma$  and (2.21) finish the proof.

We define the floating-point predecessor and successor of a real number  $r$  with  $\min\{f : f \in \mathbb{F}\} < r < \max\{f : f \in \mathbb{F}\}$  by

$$\text{pred}(r) := \max\{f \in \mathbb{F} : f < r\} \quad \& \quad \text{succ}(r) := \min\{f \in \mathbb{F} : r < f\}.$$

Using the ufp concept, the predecessor and successor of a floating-point number can be characterized as follows.

LEMMA 2.2. *Let a floating-point number  $0 \neq f \in \mathbb{F}$  be given. Then*

$$\begin{aligned} f \notin \mathbb{U} \quad \text{and} \quad |f| \neq \text{ufp}(f) &\Rightarrow \text{pred}(f) = f - 2\mathbf{eps} \cdot \text{ufp}(f) \quad \text{and} \quad f + 2\mathbf{eps} \cdot \text{ufp}(f) = \text{succ}(f) , \\ f \notin \mathbb{U} \quad \text{and} \quad f = \text{ufp}(f) &\Rightarrow \text{pred}(f) = (1 - \mathbf{eps})f \quad \text{and} \quad (1 + 2\mathbf{eps})f = \text{succ}(f) , \\ f \notin \mathbb{U} \quad \text{and} \quad f = -\text{ufp}(f) &\Rightarrow \text{pred}(f) = (1 - 2\mathbf{eps})f \quad \text{and} \quad (1 + \mathbf{eps})f = \text{succ}(f) , \\ f \in \mathbb{U} &\Rightarrow \text{pred}(f) = f - \mathbf{eta} \quad \text{and} \quad f + \mathbf{eta} = \text{succ}(f) . \end{aligned}$$

For any  $f \in \mathbb{F}$ , also in underflow,

$$(2.23) \quad \text{pred}(f) \leq f - \mathbf{eps} \cdot \text{ufp}(f) \leq f + \mathbf{eps} \cdot \text{ufp}(f) \leq \text{succ}(f) .$$

For  $f \notin \mathbb{U}$ ,

$$(2.24) \quad f - 2\mathbf{eps} \cdot \text{ufp}(f) \leq \text{pred}(f) < \text{succ}(f) \leq f + 2\mathbf{eps} \cdot \text{ufp}(f) .$$

REMARK. Note that we defined  $\mathbb{U}$  in (2.1) to contain  $\pm\frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , the smallest normalized floating-point numbers.

PROOF. For  $f \notin \mathbb{U}$  and  $|f| \neq \text{ufp}(f)$ , use  $\text{ufp}(f) < |f| < 2\text{ufp}(f)$ , and the rest is not difficult to see.  $\square$

We collect some properties when a perturbation of a floating-point number  $a$  cannot change it or, cannot go beyond the immediate neighbor of  $a$ . These will be useful in the analysis of our algorithms producing  $K$ -fold accuracy to be presented in Section 5.

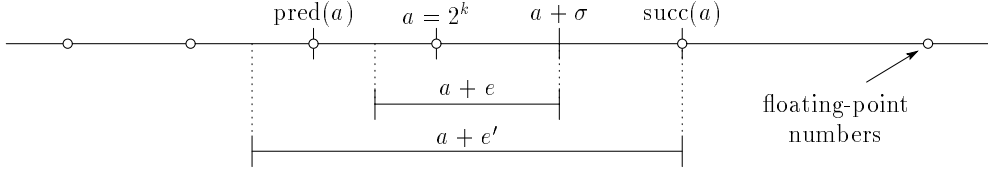


FIG. 2.2. floating-point numbers in the neighborhood of a power of 2

LEMMA 2.3. Let floating-point numbers  $a, e, e'$  be given. Define  $\mathcal{N}(a) := \{\text{pred}(a), \text{succ}(a)\}$  to be the neighbor set of  $a$ , and denote by  $N(a) \in \mathbb{F}$  one of the neighbors of  $a$ , i.e.  $N(a) \in \mathcal{N}(a)$ . Then

$$(2.25) \quad |e| < \frac{1}{2}|N(a) - a| \quad \Rightarrow \quad \text{fl}(a + e) = a ,$$

$$(2.26) \quad |a| \neq \text{ufp}(a) \quad \text{and} \quad |e| < \text{eps} \cdot \text{ufp}(a) \quad \Rightarrow \quad \text{fl}(a + e) = a .$$

Moreover,

$$(2.27) \quad \text{fl}(a + e) = a \quad \text{and} \quad |e'| < |e| + \text{eps} \cdot \text{ufp}(a) \quad \Rightarrow \quad \text{fl}(a + e') = a \quad \text{or} \quad \text{fl}(a + e') \in \mathcal{N}(a) ,$$

where  $e'e \geq 0$  in case  $\text{fl}(a + e') \neq a$ .

PROOF. The first two assertions are obvious by the definition of rounding to nearest. To see (2.27), set  $\sigma := \text{ufp}(a)$ . The critical case is that  $a$  is a power of 2, as depicted in Figure 2.2. In any case, also in the underflow range, Lemma 2.2 implies

$$\text{pred}(\text{pred}(a)) \leq a - 2\text{eps}\sigma \leq a + 2\text{eps}\sigma \leq \text{succ}(\text{succ}(a)) ,$$

and  $\text{fl}(a + e) = a$  yields

$$\frac{1}{2}(\text{pred}(a) + a) \leq a - |e| \leq a + |e| \leq \frac{1}{2}(a + \text{succ}(a)) .$$

Hence

$$\begin{aligned} \frac{1}{2}(\text{pred}(\text{pred}(a)) + \text{pred}(a)) &\leq a - |e| - \text{eps}\sigma < a - |e'| \leq \\ a + |e'| &< a + |e| + \text{eps}\sigma \leq \frac{1}{2}(\text{succ}(\text{succ}(a)) + \text{succ}(a)) . \end{aligned}$$

The lemma follows.  $\square$

The aim of this paper is to present a summation algorithm computing a faithfully rounded exact result of the sum. That means [10, 39, 9] that the computed result is definitely one of the immediate floating-point neighbors of the exact (real) result. If the exact result is a floating-point number, then this will be the result of our algorithm.

DEFINITION 2.4. A floating-point number  $f \in \mathbb{F}$  is called a faithful rounding of a real number  $r \in \mathbb{R}$  if

$$(2.28) \quad \text{pred}(f) < r < \text{succ}(f) .$$

We denote this by  $f \in \square(r)$ . For  $r \in \mathbb{F}$  this implies  $f = r$ .

For general  $r \notin \mathbb{F}$ , two floating-point numbers satisfy  $f \in \square(r)$  for given  $r \in \mathbb{R}$ , so a little accuracy is lost. The rounded to nearest result, however, has the well-known drawback that if the exact result is very close to the midpoint of two adjacent floating-point numbers, its computation may require a substantial and often not necessary additional effort. Our Algorithm 6.1 (**NearSum**) computes the rounded to nearest result. The computing time depends in this case on the exponent range of the summands.

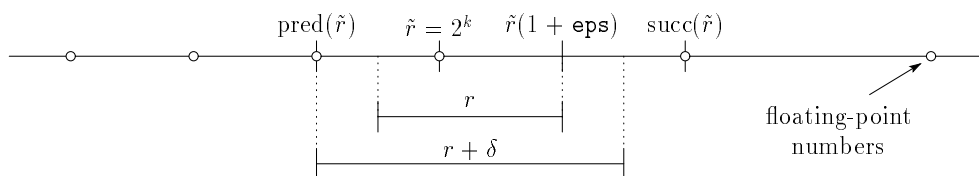


FIG. 2.3. Faithful rounding near a power of 2

In contrast, Algorithm 4.4 (**AccSum**) computes a faithfully rounded result. Its computing time is proportional to the condition number of the sum.

Suppose  $r + \delta$  is the exact result of a summation, composed of a (real) approximation  $r$  and an error term  $\delta$ . Next we establish conditions on  $\delta$  to ensure that  $\text{fl}(r)$  is a faithful rounding of  $r + \delta$ . As in (2.27), the critical case is the change of exponent at a power of 2, depicted in Figure 2.3.

**LEMMA 2.5.** *Let  $r, \delta \in \mathbb{R}$  and  $\tilde{r} := \text{fl}(r)$ . If  $\tilde{r} \notin \mathbb{U}$  suppose  $2|\delta| < \text{eps}|\tilde{r}|$ , and if  $\tilde{r} \in \mathbb{U}$  suppose  $|\delta| < \frac{1}{2}\text{eta}$ . Then  $\tilde{r} \in \square(r + \delta)$ , that means  $\tilde{r}$  is a faithful rounding of  $r + \delta$ .*

**PROOF.** According to Definition 2.4 we have to prove  $\text{pred}(\tilde{r}) < r + \delta < \text{succ}(\tilde{r})$ . If  $\tilde{r} \in \mathbb{U}$ , then  $|\tilde{r} - r| \leq \frac{1}{2}\text{eta}$  by (2.19), so Lemma 2.2 yields

$$\text{pred}(\tilde{r}) = \tilde{r} - \text{eta} < \tilde{r} - |\tilde{r} - r| + \delta \leq r + \delta \leq \tilde{r} + |\tilde{r} - r| + \delta < \tilde{r} + \text{eta} = \text{succ}(\tilde{r})$$

and finishes this part. It remains to treat the case  $\tilde{r} \notin \mathbb{U}$ .

Then  $\text{ufp}(\tilde{r}) \leq |\tilde{r}| < 2\text{ufp}(\tilde{r})$  by (2.9), so  $|\delta| < \text{eps} \cdot \text{ufp}(\tilde{r})$ . Suppose  $r \leq \tilde{r}$ . Then rounding to nearest implies

$$0 \leq \tilde{r} - r \leq \frac{1}{2}(\tilde{r} - \text{pred}(\tilde{r})) \quad \text{and} \quad |\delta| < \frac{1}{2}(\tilde{r} - \text{pred}(\tilde{r})),$$

where the latter follows directly from Lemma 2.2 if  $|\tilde{r}|$  is not a power of 2, and otherwise by  $2|\delta| < \text{eps}|\tilde{r}| = \text{eps} \cdot \text{ufp}(\tilde{r}) = \tilde{r} - \text{pred}(\tilde{r})$ . Hence (2.23) yields

$$\text{pred}(\tilde{r}) = \tilde{r} - (\tilde{r} - \text{pred}(\tilde{r})) < \tilde{r} - (\tilde{r} - r) - |\delta| \leq r + \delta \leq \tilde{r} + \delta < \tilde{r} + \text{eps} \cdot \text{ufp}(\tilde{r}) \leq \text{succ}(\tilde{r}).$$

The case  $r > \tilde{r}$  follows similarly. □

A faithfully rounded result satisfies some weak ordering properties. For  $f, f_1, f_2 \in \mathbb{F}$ ,  $r \in \mathbb{R}$  and  $f \in \square(r)$ , i.e.  $f$  is a faithful rounding of  $r$ , one verifies

$$(2.29) \quad \begin{aligned} f_1 < f < f_2 &\Rightarrow f_1 < r < f_2 \\ f_1 < r < f_2 &\Rightarrow f_1 \leq f \leq f_2. \end{aligned}$$

The following lemma shows that if the sum  $s \in \mathbb{R}$  of some floating-point numbers is in the underflow range, a faithful rounding  $f$  of  $s$  is strong enough to imply  $f = s$ .

**LEMMA 2.6.** *Let  $r \in \text{eta}\mathbb{Z}$  be given, and let  $f \in \mathbb{F}$  be a faithful rounding of  $r$ . Then*

$$(2.30) \quad |r| \leq \frac{1}{2}\text{eps}^{-1}\text{eta} \Rightarrow f \in \mathbb{U} \Rightarrow r = f.$$

Therefore especially  $f = 0$  iff  $r = 0$  and

$$\text{sign}(r) = \text{sign}(f).$$

For  $r \neq 0$  it holds

$$(2.31) \quad |r - f| < 2\text{eps}|r| \quad \text{and} \quad |r - f| < 2\text{eps} \cdot \text{ufp}(f).$$

PROOF. If  $|r| = \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , then faithful rounding implies  $f = r \in \mathbb{U}$ . If  $|r| < \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , then  $f \in \mathbb{U}$  by (2.29), so that  $f \pm \mathbf{eta}$  are the neighbors of  $f$  by Lemma 2.2. Hence (2.28) yields  $|r - f| < \mathbf{eta}$ , and  $f \in \mathbb{F} \subseteq \mathbf{eta}\mathbb{Z}$  and  $r \in \mathbf{eta}\mathbb{Z}$  imply  $r = f$ . This proves (2.30).

To prove (2.31) we may assume  $f \notin \mathbb{U}$  and without loss of generality  $f > 0$ . For  $f$  not being a power of 2 we have  $\mathbf{ufp}(f) < |r|$ , so Lemma 2.2 and (2.28) imply

$$|r - f| < 2\mathbf{eps} \cdot \mathbf{ufp}(f) \leq 2\mathbf{eps}|r| .$$

This is also true for  $f$  being a power of 2 and  $f \leq r < \mathbf{succ}(f)$ , and for  $(1 - \mathbf{eps})f = \mathbf{pred}(f) < r \leq f$  we have

$$|r - f| < \mathbf{eps} \cdot \mathbf{ufp}(f) = \mathbf{eps} \cdot f < (1 - \mathbf{eps})^{-1}\mathbf{eps} \cdot r < 2\mathbf{eps}|r| .$$

The lemma is proved. □

Definition 2.4 and Lemma 2.5 are formulated for general  $r \in \mathbb{R}$ . For our main application, the approximation of the sum  $s = \sum a_i$  of floating-point numbers,  $a_i \in \mathbf{eta}\mathbb{Z}$  implies  $s \in \mathbf{eta}\mathbb{Z}$  and Lemma 2.6 is applicable.

As has been noted in (2.4), the error of a floating-point addition is always a floating-point number. Fortunately, rather than Algorithm 2.1 (**TwoSum**) we can use in our applications the following faster algorithm due to Dekker [10]. Again, the computation is very efficient because only standard floating-point addition and subtraction is used.

ALGORITHM 2.7. *Compensated summation of two floating-point numbers.*

```
function [x, y] = FastTwoSum(a, b)
    x = fl(a + b)
    q = fl(x - a)
    y = fl(b - q)
```

In Dekker's original algorithm,  $y$  is computed by  $y = \mathbf{fl}((a - x) + b)$ , which is equivalent to the last statement in Algorithm 2.7 because  $\mathbb{F} = -\mathbb{F}$  and  $\mathbf{fl}(-r) = -\mathbf{fl}(r)$  for  $r \in \mathbb{R}$ . For floating-point arithmetic with rounding to nearest and base 2, e.g. IEEE 754 arithmetic, Dekker [10] showed in 1971 that the correction is *exact* if the input is ordered by magnitude, that is

$$(2.32) \quad x + y = a + b$$

provided  $|a| \geq |b|$ . In [36] we showed that the obvious way to get rid of this assumption is suboptimal on today's computers because a branch slows down computation significantly (see also Section 9).

Algorithm 2.7 (**FastTwoSum**) is an error-free transformation of the pair of floating-point numbers  $(a, b)$  into a pair  $(x, y)$ . Our first algorithm 4.4 (**AccSum**) to be presented can also be viewed as an error-free transformation of a vector  $p$  into floating-point numbers  $\tau_1, \tau_2$  and a vector  $p'$  such that  $\sum p_i = \tau_1 + \tau_2 + \sum p'_i$ , and  $\mathbf{res} := \mathbf{fl}(\tau_1 + (\tau_2 + \sum p'_i))$  is the faithfully rounded sum  $\sum p_i$ . To prove this we need to refine the analysis of Algorithm 2.7 by weakening the assumption  $|a| \geq |b|$ . For completeness we first prove Dekker's result [10].

LEMMA 2.8. *Let  $a, b$  be floating-point numbers and  $|a| \geq |b|$ . Let  $x, y$  be the results produced by Algorithm 2.7 (**FastTwoSum**) applied to  $a, b$ . Then*

$$(2.33) \quad x + y = a + b, \quad x = \mathbf{fl}(a + b) \quad \text{and} \quad |y| \leq \mathbf{eps} \cdot \mathbf{ufp}(a + b) \leq \mathbf{eps} \cdot \mathbf{ufp}(x) .$$

Furthermore,

$$(2.34) \quad q = \mathbf{fl}(x - a) = x - a \quad \text{and} \quad y = \mathbf{fl}(b - q) = b - q ,$$

```

for  $i = 1 : n$ 
   $q = \text{fix}(p_i/b)$   % fix: Round to integer towards zero
   $t_j = t_j + q$ 
   $p_i = p_i - q * b$ 
end for

```

FIG. 3.1. Inner loop of the summation by Zielke and Drygalla

that is the floating-point subtractions  $x - a$  and  $b - q$  are exact.

PROOF. We first show  $\text{fl}(x - a) = x - a$ . Without loss of generality we can assume  $a \geq 0$ . If  $0 \leq b \leq a$ , then  $a + b \leq 2a$  and  $x = \text{fl}(a + b) \leq 2a$  by (2.6), so  $\frac{1}{2}x \leq a \leq x$  and (2.22) imply  $\text{fl}(x - a) = -\text{fl}(a - x) = x - a$ . If  $\frac{1}{2}a \leq -b \leq a$ , then  $-x = \text{fl}(-b - a) = -b - a$  again by (2.22), and  $\text{fl}(x - a) = \text{fl}(b) = b = x - a$ . If finally  $0 \leq -b < \frac{1}{2}a$ , then  $\frac{1}{2}a < a + b \leq a$  and  $\frac{1}{2}a \leq x \leq a$  by (2.6), so  $\text{fl}(x - a) = x - a = q$  again by (2.22). This proves  $q = x - a$ .

Define  $\delta := x - (a + b) = q - b$ . Then  $\delta \in \mathbb{F}$  by (2.4) and  $y = \text{fl}(b - q) = \text{fl}(-\delta) = -\delta = b - q$ , proving (2.34) and the first part of (2.33). The second part follows by  $|y| = |\delta|$  and (2.20).  $\square$

Next we weaken the assumption  $|a| \geq |b|$ .

LEMMA 2.9. *Let  $a, b$  be floating-point numbers with  $\text{ufp}(b) \leq \text{ufp}(a)$ . Let  $x, y$  be the results produced by Algorithm 2.7 (FastTwoSum) applied to  $a, b$ . Then the assertions (2.33) and (2.34) of Lemma 2.8 are still valid.*

PROOF. Using Lemma 2.8 we can assume  $|a| < |b|$ . This implies  $\text{ufp}(b) \leq \text{ufp}(a) \leq |a| < |b| < 2\text{ufp}(b)$ , so  $\text{ufp}(b) = \text{ufp}(a) =: \sigma/2$ , and  $a, b \in \text{eps}\sigma\mathbb{Z}$  by (2.13). Hence  $|a + b| < 2\text{ufp}(a) + 2\text{ufp}(b) = 2\sigma$ , so  $\text{ufp}(a + b) \leq \sigma$ . Furthermore, (2.15) and (2.20) yield  $\delta := x - (a + b) \in \text{eps}\sigma\mathbb{Z}$  and  $|\delta| \leq \text{eps} \cdot \text{ufp}(a + b) \leq \text{eps}\sigma$ . That means, the binary representation of  $\delta$  has at most one nonzero bit. Furthermore,  $|b| < \sigma$  and (2.23) imply  $|b| \leq \text{pred}(\sigma) \leq \sigma - \text{eps}\sigma$ , also in the presence of underflow, and therefore  $|b + \delta| \leq \sigma$ . So (2.21) yields  $x - a = b + \delta = \text{fl}(b + \delta) = \text{fl}(x - a) = q$ , and  $y = \text{fl}(b - q) = \text{fl}(-\delta) = -\delta = b - q$ , proving (2.34). Hence  $x + y = x - \delta = a + b$ , and the estimation on  $|y| = |\delta|$  follows by (2.20). This finishes the proof.  $\square$

Lemma 2.9 may also offer possibilities for summation algorithms based on sorting: To apply FastTwoSum it suffices to “sort by exponent”, which has complexity  $\mathcal{O}(n)$ .

Lemma 2.9 shows that  $|a|$  may be less than  $|b|$  without jeopardizing the assertions of Lemma 2.8 as long as the leading bit of  $a$  is not less than the one of  $b$ . The next lemma shows that also this assumption may be dropped as long as there are not too many trailing nonzero bits of  $a$ .

LEMMA 2.10. *Let  $a, b$  be floating-point numbers and  $\sigma = 2^k$  for  $k \in \mathbb{Z}$ . Suppose  $a, b \in \text{eps}\mathbb{Z}$  and  $|b| < \sigma$ . Let  $x, y$  be the results produced by Algorithm 2.7 (FastTwoSum) applied to  $a, b$ . Then the assertions (2.33) and (2.34) of Lemma 2.8 are still valid.*

PROOF. Let  $\text{fl}(a + b) = a + b + \delta$ . If  $|a| \geq \frac{1}{2}\sigma$ , we can use Lemma 2.9. Otherwise,  $|a + b| < \frac{3}{2}\sigma$ , so (2.15) and (2.20) imply  $|\delta| \leq \text{eps}\sigma$ . In fact, either  $|\delta| = \text{eps}\sigma$  or  $\delta = 0$ . Hence  $|x - a| = |b + \delta| \leq \text{pred}(\sigma) + \text{eps}\sigma \leq \sigma$ , so (2.21) yields  $q = \text{fl}(x - a) = x - a$ . Now we can proceed as in the proof of Lemma 2.9.  $\square$

**3. Extraction of high order parts.** The main idea of the method by Zielke and Drygalla [45] to compute the sum of floating-point numbers  $p_i, 1 \leq i \leq n$ , is as follows. First they compute  $k \in \mathbb{N}$  such that  $\max |p_i| < 2^k$ , and some  $M$  such that  $n < 2^M$ . Then, for  $\text{eps} = 2^{-53}$ , they extract the “bits” from  $k$  down to  $k - (53 - M) + 1$  of  $p_i$  into  $q_i$  and add the  $q_i$  into some  $t_1$ . By the choice of  $M$  this addition is exact. They

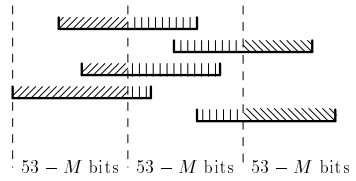


FIG. 3.2. Zielke and Drygalla’s scheme

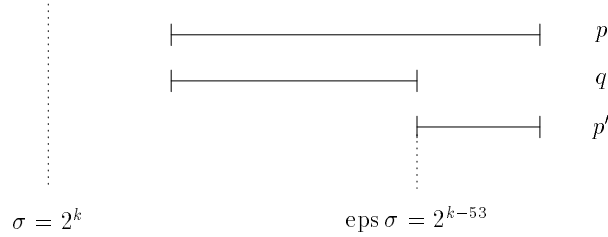


FIG. 3.3. Extract Scalar: error-free transformation  $p = q + p'$

then continue by extracting the “bits” from  $k - (53 - M)$  downto  $k - 2(53 - M) + 1$  and summing them into  $t_2$  and so forth until all bits are processed. For given  $M$  and  $k$  and assuming no underflow the inner loop looks as in Figure 3.1. Since  $q$  can be overwritten, no index is necessary. Here  $b$  is an appropriate power of 2, which is updated before entering the for-loop again. By the principle of the approach, the sum  $t_j$  of one “chunk” of  $53 - M$  bits is exact. After finishing, the  $t_j$ , appropriately shifted, are added.

The method is illustrated in Figure 3.2. The authors also mention a faster, though less accurate method, to extract only the leading  $53 - M$  bits and adding the rest accepting accumulating rounding errors. However, for neither method an error analysis is given.

The inner loop of the method requires to split a floating-point number  $p \in \mathbb{F}$  according to Figure 3.3, which is done by chopping (fix) by Zielke and Drygalla. There are other possibilities, for example to round from floating-point to integer (rather than chopping), or the assignment of a floating-point number to a long integer (if supported by the hardware in use). Also direct manipulation by accessing mantissa and exponent is possible. However, all these methods slow down the extraction significantly, often by an order of magnitude and more compared to our following Algorithm 3.1 (`ExtractScalar`). We will show corresponding performance data in Section 9.

Our approach follows a similar scheme as in Figure 3.2. However, we carefully estimate how many bits of the result have to be extracted to guarantee a faithful rounding of the result. For the splitting as depicted in Figure 3.3, neither the high order part  $q$  and low order part  $p'$  need to match bitwise with the original  $p$ , nor must  $q$  and  $p'$  have the same sign; only the error-freeness of the transformation  $p = q + p'$  is mandatory. This is achieved by the following fast algorithm, where  $\sigma$  denotes a power of 2 not less than  $|p|$ .

ALGORITHM 3.1. *Error-free transformation extracting high order part.*

$$\begin{aligned} \text{function } [q, p'] &= \text{ExtractScalar}(\sigma, p) \\ q &= \text{fl}((\sigma + p) - \sigma) \\ p' &= \text{fl}(p - q) \end{aligned}$$

There is an important difference to the splitting in Algorithm 1.2 (`Split`). There, a floating-point number is split into two parts relative to its exponent, and both the high and the low part has at most 26 significant

bits. In `ExtractScalar` a floating-point number is split relative to  $\sigma$ , a fixed power of 2. The higher and the lower part of the splitting may have between 0 and 53 significant bits, depending on  $\sigma$ .

This clever way of splitting in Algorithm 3.1 (`ExtractScalar`) is due to the second author. The performance of splitting is crucial since it is in the inner loops of our algorithms.

We think this method is known, but the only reference we found for this is [16]. However, we do not know of any analysis of Algorithm 3.1, so we develop it in the following lemma.

**LEMMA 3.2.** *Let  $q$  and  $p'$  be the results of Algorithm 3.1 (`ExtractScalar`) applied to floating-point numbers  $\sigma$  and  $p$ . Assume  $\sigma = 2^k \in \mathbb{F}$  for some  $k \in \mathbb{Z}$ , assume  $|p| \leq 2^{-M}\sigma$  for some  $0 \leq M \in \mathbb{N}$ . Then*

$$(3.1) \quad p = q + p' , \quad |p'| \leq \mathbf{eps}\sigma , \quad |q| \leq 2^{-M}\sigma \quad \text{and} \quad q \in \mathbf{eps}\sigma\mathbb{Z} .$$

**PROOF.** We first note that `ExtractScalar`( $\sigma, p$ ) performs exactly the same operations in the same order as  $[x, y] = \mathbf{FastTwoSum}(\sigma, p)$ , so  $|p| \leq \sigma$  and Lemma 2.8 imply  $p' = p - q$ . If  $|p| = \sigma$ , then  $p' = 0$ , otherwise (2.33) implies  $|p'| \leq \mathbf{eps} \cdot \mathbf{ufp}(\sigma + p) \leq \mathbf{eps}\sigma$ . Furthermore,  $q \in \mathbf{eps}\sigma\mathbb{Z}$  follows by (2.17).

Finally we prove  $|q| \leq 2^{-M}\sigma$ . First, suppose  $\sigma + \mathbf{sign}(p)2^{-M}\sigma$  is a floating-point number. Then  $\sigma + p$  is in the interval with floating-point endpoints  $\sigma$  and  $\sigma + \mathbf{sign}(p)2^{-M}\sigma$ , so (2.6) implies that  $x := \mathbf{fl}(\sigma + p)$  is in that interval and  $|q| = |x - \sigma| \leq 2^{-M}\sigma$  follows. Second, suppose  $\sigma + \mathbf{sign}(p)2^{-M}\sigma$  is not a floating-point number. Then  $\mathbf{fl}(\sigma + \mathbf{sign}(p)2^{-M}\sigma) = \sigma$  because  $\mathbf{sign}(p)2^{-M}\sigma$  is a power of 2 and rounding to nearest is tie to even, so monotonicity of the rounding implies  $\mathbf{fl}(\sigma + p) = \sigma$  and  $q = 0$ .  $\square$

Following we adapt Algorithm 3.1 (`ExtractScalar`) to the error-free transformation of an entire vector. For this case we prove that the high order parts can be summed up without error. For better readability and analysis the extracted parts are stored in a vector  $q_i$ . In a practical implementation, the vector  $q$  is not necessary but only its sum.

**ALGORITHM 3.3.** *Error-free vector transformation extracting high order part.*

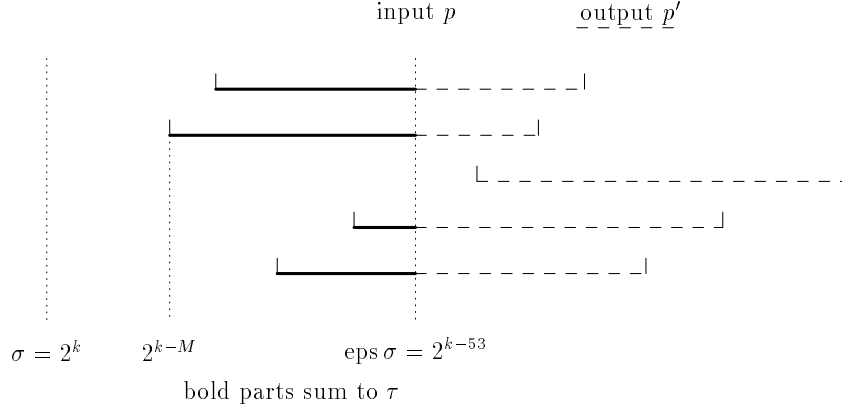
```
function  $[\tau, p'] = \mathbf{ExtractVector}(\sigma, p)$ 
   $\tau = 0$ 
  for  $i = 1 : n$ 
     $[q_i, p'_i] = \mathbf{ExtractScalar}(\sigma, p_i)$ 
     $\tau = \mathbf{fl}(\tau + q_i)$ 
  end for
```

Algorithm 3.3 proceeds as depicted in Figure 3.4. Note that the loop can be executed in parallel.

Note again that the low order parts, which are collected in  $p'$ , neither need to be bitwise identical to those of  $p$  nor do they need to have the same sign. The important property we are after is that the transformation is performed without error, i.e.  $\sum p_i = \tau + \sum p'_i$ , and that  $|p'_i|$  stays below  $\mathbf{eps}\sigma$ . The validity of the algorithm is demonstrated by the following theorem.

**THEOREM 3.4.** *Let  $\tau$  and  $p'$  be the results of Algorithm 3.3 (`ExtractVector`) applied to  $\sigma \in \mathbb{F}$  and a vector of floating-point numbers  $p_i, 1 \leq i \leq n$ . Assume  $\sigma = 2^k \in \mathbb{F}$  for some  $k \in \mathbb{Z}$ ,  $n + 2 \leq 2^M$  for some  $M \in \mathbb{N}$  and  $|p_i| \leq 2^{-M}\sigma$  for all  $i$ . Then*

$$(3.2) \quad \sum_{i=1}^n p_i = \tau + \sum_{i=1}^n p'_i , \quad \max |p'_i| \leq \mathbf{eps}\sigma , \quad |\tau| < \sigma \quad \text{and} \quad \tau \in \mathbf{eps}\sigma\mathbb{Z} .$$


 FIG. 3.4. *Extract vector: error-free transformation*  $\sum p_i = \tau + \sum p'_i$ 

If  $2^{2M-1}\mathbf{eps} \leq 1$  and  $\sigma \notin \mathbb{U}$ , then

$$(3.3) \quad |\mathbf{fl}(\tau + T)| < \sigma \quad \text{for} \quad T := \mathbf{fl}\left(\sum_{i=1}^n p'_i\right).$$

REMARK. Note that for (3.2) we have no assumption on the size of  $M$ . However, for extremely large  $M$  with  $2^M\mathbf{eps} \geq 1$ ,  $|p_i| \leq 2^{-M}\sigma$  implies that the extracted parts summed up into  $\tau$  consist of at most one bit, so not much is gained.

PROOF OF THEOREM 3.4. Lemma 3.2 implies  $p_i = q_i + p'_i$ ,  $|p'_i| \leq \mathbf{eps}\sigma$ ,  $|q_i| \leq 2^{-M}\sigma$  and  $q_i \in \mathbf{eps}\sigma\mathbb{Z}$  for all  $i \in \{1, \dots, n\}$ . So  $\tau \in \mathbf{eps}\sigma\mathbb{Z}$ , and  $\sum |q_i| \leq n2^{-M}\sigma < \sigma$  and (2.21) imply  $\tau = \mathbf{fl}(\sum q_i) = \sum q_i$  and  $|\tau| \leq n2^{-M}\sigma < \sigma$ . This proves (3.2). By (3.2) and (2.16) it follows  $|T| \leq n\mathbf{eps}\sigma$ , so Lemma 2.2 implies

$$\begin{aligned} |\tau + T| &\leq (n2^{-M} + n\mathbf{eps})\sigma \leq (2^M - 2)(2^{-M} + \mathbf{eps})\sigma \\ &= (1 - 2^{-M+1}(1 - 2^{2M-1}\mathbf{eps}) - 2\mathbf{eps})\sigma < (1 - \mathbf{eps})\sigma \\ &= \mathbf{pred}(\sigma) \end{aligned}$$

because  $\sigma \notin \mathbb{U}$ , so that (2.6) implies  $|\mathbf{fl}(\tau + T)| = \mathbf{fl}(|\tau + T|) \leq \mathbf{pred}(\sigma) < \sigma$ .  $\square$

Note that (3.2) requires no assumption on  $\sigma$  other than being a power of 2 and bounding  $2^M|p_i|$ ;  $\sigma$  may well be in the underflow range.

To apply Theorem 3.4 the best (smallest) values for  $M$  and  $\sigma$  are  $\lceil \log_2(n+2) \rceil$  and  $2^{M+\lceil \log_2(\max |p_i|) \rceil}$ , respectively. To avoid the use of standard functions, these can be calculated by the following algorithm.

ALGORITHM 3.5. *Computation of*  $2^{\lceil \log_2 |p| \rceil}$  *for*  $p \neq 0$ .

```

function L = NextPowerTwo(p)
    q = eps-1p
    L = fl(|(q + p) - q|)
    if L = 0
        L = |p|
    end if
    
```

THEOREM 3.6. *Let*  $L$  *be the result of Algorithm 3.5 (NextPowerTwo) applied to a nonzero floating-point number*  $p$ . *If no overflow occurs, then*  $L = 2^{\lceil \log_2 |p| \rceil}$ .

REMARK. For simplicity we skipped the obvious check for large input number  $p$  to avoid overflow in the computation of  $q$ . However, we will show that  $L = 2^{\lceil \log_2 |p| \rceil}$  is satisfied in the presence of underflow.

PROOF. First assume  $|p| = 2^k$  for some  $k \in \mathbb{Z}$ . Then the rounding tie to even implies  $\text{fl}(q + p) = \text{fl}(q(1 + \text{eps})) = q$ , so that  $\text{fl}(|(q + p) - q|) = 0$ , and  $L = |p| = 2^k = 2^{\lceil \log_2 |p| \rceil}$  for the final result  $L$ .

So we may assume that  $p$  is not a power of 2, and without loss of generality we assume  $p > 0$ . Then

$$\text{ufp}(p) < p < 2\text{ufp}(p),$$

and we have to show  $L = 2\text{ufp}(p)$ . Define  $x := \text{fl}(q + p)$ . By Lemma 2.8 the computation of  $\text{fl}(x - q)$  causes no rounding error, so that  $L = \text{fl}(q + p) - q$ . By definition,  $\text{ufp}(q) = \text{eps}^{-1}\text{ufp}(p) < \text{eps}^{-1}p = q$ , so that  $q \notin \mathbb{U}$  and Lemma 2.2 imply  $\text{succ}(q) = q + 2\text{eps} \cdot \text{ufp}(q)$ . That means  $q + \text{eps} \cdot \text{ufp}(q)$  is the midpoint of  $q$  and  $\text{succ}(q)$ . Hence rounding to nearest and

$$q + \text{eps} \cdot \text{ufp}(q) < q + \text{eps} \cdot q = q + p < q + 2\text{ufp}(p) = \text{succ}(q)$$

implies  $\text{fl}(q + p) = \text{succ}(q)$ , so that  $L = \text{fl}(q + p) - q = 2\text{eps} \cdot \text{ufp}(q) = 2\text{ufp}(p)$ . The theorem is proved.  $\square$

**4. Algorithms and analysis.** To ease analysis, we formulate our summation algorithms with superscripts to variables to identify the different stages. Of course in the actual implementation especially vectors are overwritten.

We first transform an input vector plus offset into two leading and a remaining part. The following formulation is aimed on readability rather than efficiency. In this first step we especially avoid a check for zero sum.

ALGORITHM 4.1. *Preliminary version of transformation of vector  $p^{(0)}$  plus offset  $\varrho$ .*

```

function  $[\tau_1, \tau_2, p^{(m)}, \sigma] = \text{Transform}(p^{(0)}, \varrho)$ 
   $\mu = \max(|p_i^{(0)}|)$ 
  if  $\mu = 0$ ,  $\tau_1 = \varrho$ ,  $\tau_2 = p^{(m)} = \sigma = 0$ , return, end if
   $M = \lceil \log_2 (\text{length}(p^{(0)}) + 2) \rceil$ 
   $\sigma_0 = 2^{M + \lceil \log_2(\mu) \rceil}$ 
   $t^{(0)} = \varrho$ ,  $m = 0$ 
  repeat
     $m = m + 1$ 
     $[\tau^{(m)}, p^{(m)}] = \text{ExtractVector}(\sigma_{m-1}, p^{(m-1)})$ 
     $t^{(m)} = \text{fl}(t^{(m-1)} + \tau^{(m)})$ 
     $\sigma_m = \text{fl}(2^M \text{eps} \sigma_{m-1})$ 
  until  $\sigma_{m-1} \leq \frac{1}{2} \text{eps}^{-1} \text{eta}$  or  $|t^{(m)}| \geq \text{fl}(2^{2M+1} \text{eps} \sigma_{m-1})$ 
   $\sigma = \sigma_{m-1}$ 
   $[\tau_1, \tau_2] = \text{FastTwoSum}(t^{(m-1)}, \tau^{(m)})$ 

```

REMARK 1. The output parameter  $\sigma$  is not necessary in the following applications of **Transform** but added for clarity in the forthcoming proofs.

REMARK 2. For clarity we also use for the moment the logarithm in the computation of  $M$  and  $\sigma_0$ . Later this will be replaced by Algorithm 3.5 (**NextPowerTwo**) (see the Matlab code given in the Appendix).

LEMMA 4.2. *Let  $\tau_1, \tau_2, p^{(m)}, \sigma$  be the results of Algorithm 4.1 (**Transform**) applied to a nonzero vector of floating-point numbers  $p_i^{(0)}$ ,  $1 \leq i \leq n$ , and  $\varrho \in \mathbb{F}$ . Assume that  $2(n+2)^2 \text{eps} \leq 1$  and  $\varrho \in \text{eps} \sigma_0 \mathbb{Z}$  is satisfied for  $M := \lceil \log_2(n+2) \rceil$ ,  $\mu := \max_i |p_i|$  and  $\sigma_0 = 2^{M + \lceil \log_2 \mu \rceil}$ . Denote  $s := \sum_{i=1}^n p_i^{(0)}$ .*

Then Algorithm 4.1 will stop, and

$$(4.1) \quad s + \varrho = t^{(m-1)} + \tau^{(m)} + \sum_{i=1}^n p_i^{(m)}$$

$$(4.2) \quad \max |p_i^{(m)}| \leq \mathbf{eps}\sigma_{m-1}, \quad |\tau^{(m)}| < \sigma_{m-1} \quad \text{and} \quad t^{(m-1)}, \tau^{(m)} \in \mathbf{eps}\sigma_{m-1}\mathbb{Z}$$

is true for all  $m$  between 1 and its final value. Moreover,

$$(4.3) \quad \tau_1 + \tau_2 = t^{(m-1)} + \tau^{(m)}, \quad \tau_1 = \mathbf{fl}(\tau_1 + \tau_2) = \mathbf{fl}(t^{(m-1)} + \tau^{(m)}) = t^{(m)}$$

is satisfied for the final value of  $m$ . If  $\sigma_{m-1} > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  is satisfied for the final value of  $m$ , then

$$(4.4) \quad \mathbf{ufp}(\tau_1) \geq 2^{2M+1}\mathbf{eps}\sigma_{m-1}.$$

REMARK. Note that the computation of  $\sigma_m$  may be afflicted with a rounding error if  $\sigma_{m-1}$  is in the underflow range  $\mathbb{U}$ . However, we will see that this cannot do any harm. The computation of  $\sigma = \sigma_{m-1}$  can never be afflicted with rounding errors.

PROOF OF LEMMA 4.2. Algorithm **Transform** will stop because  $\sigma_m$  is decreased in each loop. We proceed by induction to prove (4.1) and (4.2). Theorem 3.6 and the initialization in Algorithm 4.1 imply  $M = \lceil \log_2(n+2) \rceil$  and  $\max |p_i^{(0)}| = \mu \leq 2^{\lceil \log_2(\mu) \rceil} = 2^{-M}\sigma_0$ , so that the assumptions of Theorem 3.4 are satisfied for  $\sigma_0$  and  $p^{(0)}$  as input to **ExtractVector**. Note this is also true if  $\sigma_0 \in \mathbb{U}$ . This yields (4.2) for  $m = 1$ , where  $t^{(0)} = \varrho \in \mathbf{eps}\sigma_0\mathbb{Z}$  is assured by the assumptions of the lemma. Furthermore,  $s = \tau^{(1)} + \sum p_i^{(1)}$ , and (4.1) is also proved for  $m = 1$ .

Now assume the repeat-until loop has been executed, denote by  $m$  the current value (immediately before the “until”-statement), and assume that (4.1) and (4.2) are true for the previous index  $m-1$ . The previous “until”-condition in Algorithm 4.1 implies

$$\sigma_{m-2} \geq \mathbf{eps}^{-1}\mathbf{eta}$$

because  $\sigma_{m-2}$  is a power of 2, hence no rounding error has occurred in the previous computation of  $2^{2M+1}\mathbf{eps}\sigma_{m-2}$ . So the induction hypothesis and the “until”-condition yield

$$(4.5) \quad |t^{(m-1)}| = |\mathbf{fl}(t^{(m-2)} + \tau^{(m-1)})| < 2^{2M+1}\mathbf{eps}\sigma_{m-2} \leq \sigma_{m-2},$$

the latter by the assumptions on  $n$  and  $M$ . By induction hypothesis,  $t^{(m-2)}, \tau^{(m-1)} \in \mathbf{eps}\sigma_{m-2}\mathbb{Z}$ , so that (2.21) implies

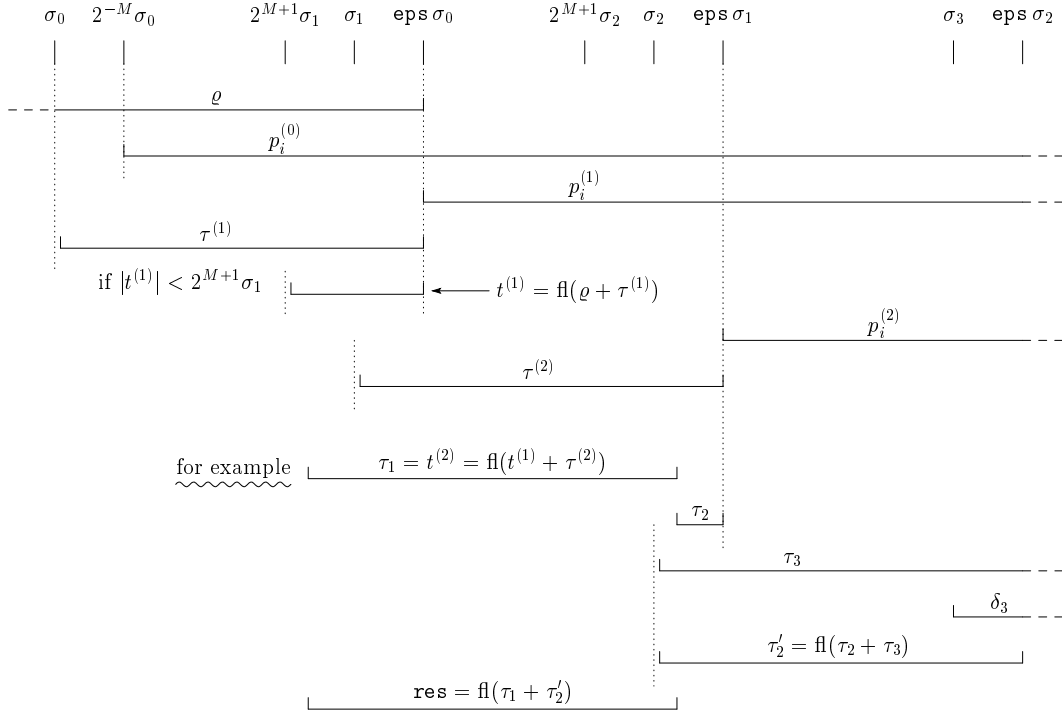
$$(4.6) \quad t^{(m-1)} = \mathbf{fl}(t^{(m-2)} + \tau^{(m-1)}) = t^{(m-2)} + \tau^{(m-1)},$$

and the induction hypothesis on (4.1) yields

$$s + \varrho = t^{(m-2)} + \tau^{(m-1)} + \sum_{i=1}^n p_i^{(m-1)} = t^{(m-1)} + \sum_{i=1}^n p_i^{(m-1)}.$$

By (4.2) we know  $\max |p_i^{(m-1)}| \leq \mathbf{eps}\sigma_{m-2} = 2^{-M}\sigma_{m-1}$ . Hence Theorem 3.4 is applicable and shows  $\sum p_i^{(m-1)} = \tau^{(m)} + \sum p_i^{(m)}$ , and therefore (4.1) for index  $m$ . It also shows (4.2), where  $t^{(m-1)} \in \mathbf{eps}\sigma_{m-1}\mathbb{Z}$  follows by  $t^{(m-2)}, \tau^{(m-1)} \in \mathbf{eps}\sigma_{m-2}\mathbb{Z}$ , (4.6) and (2.11). We proved (4.1) and (4.2). Therefore, for the last line in Algorithm **Transform** the assumptions of Lemma 2.10 are satisfied and (4.3) follows.

If  $\sigma_{m-1} > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  is satisfied for the final value of  $m$ , then  $\sigma_{m-1} \geq \mathbf{eps}^{-1}\mathbf{eta}$  because  $\sigma_{m-1}$  is a power of 2. Therefore,  $\mathbf{fl}(2^{2M+1}\mathbf{eps}\sigma_{m-1}) = 2^{2M+1}\mathbf{eps}\sigma_{m-1}$ , and the “until”-condition and (4.3) yield  $|\tau_1| = |t^{(m)}| \geq 2^{2M+1}\mathbf{eps}\sigma_{m-1}$ , which implies (4.4). The lemma is proved.  $\square$

FIG. 4.1. Outline of faithful rounding for  $m = 2$ .

The proof of Lemma 4.2 makes also clear why the assumption  $\rho \in \mathbf{eps}\sigma_0\mathbb{Z}$  is necessary. If Algorithm **Transform** does not stop for  $m = 1$ , then we showed that no rounding error occurs in the computation of  $t^{(1)} = \mathbf{fl}(t^{(0)} + \tau^{(1)})$ . Without the assumption on  $\rho = t^{(0)}$  this need not be true since  $\rho$  could be less than  $\mathbf{eps}|\tau^{(1)}|$ .

The case  $s = 0$  is far from being treated optimal. In this case the preliminary version of Algorithm 4.1 *always* iterates until  $\sigma_{m-1} \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , and each time a vector  $p$  is extracted which may long consist only of zero components. The case  $s = 0$  is not that rare, for example when checking geometrical predicates. We will improve on that later.

Next we will show how to compute a faithfully rounded result. In Figure 4.1 we sketch the possible spread of bits of the individual variables in Algorithm 4.1 (**Transform**) for a final value  $m = 2$ . In the middle of the figure we define a possible  $\tau_1 = t^{(m)}$ , since the following quantities  $\tau_2$  etc. depend on that. Note this is a picture for small  $M$ , i.e. small dimension  $n$ . For larger  $n$ , up to almost  $\sqrt{\frac{1}{2}\mathbf{eps}^{-1}}$ , the error  $\delta_3$  in the computation of  $\tau_3 = \mathbf{fl}(\sum p_i^{(m)})$  can be quite significant, though still just not too large to jeopardize faithful rounding.

**LEMMA 4.3.** *Let  $p$  be a nonzero vector of  $n$  floating-point numbers and  $\rho \in \mathbb{F}$ . Let  $\mathbf{res}$  be computed as follows:*

$$\begin{aligned} [\tau_1, \tau_2, p', \sigma] &= \mathbf{Transform}(p, \rho) \\ \mathbf{res} &= \mathbf{fl}(\tau_1 + (\tau_2 + (\sum_{i=1}^n p'_i))) \end{aligned}$$

*Assume that  $2(n+2)^2\mathbf{eps} \leq 1$ , and assume  $\rho \in \mathbf{eps}\sigma_0\mathbb{Z}$  is satisfied for  $M := \lceil \log_2(n+2) \rceil$ ,  $\mu := \max_i |p_i|$  and  $\sigma_0 = 2^{M+\lceil \log_2 \mu \rceil}$ .*

*Then  $\mathbf{res}$  is a faithful rounding of  $\sum_{i=1}^n p_i + \rho =: s + \rho$ . Moreover,*

$$(4.7) \quad s + \rho = \tau_1 + \tau_2 + \sum_{i=1}^n p'_i \quad \text{and} \quad \max |p'_i| \leq \mathbf{eps}\sigma,$$

$$(4.8) \quad \text{fl}(\tau_1 + \tau_2) = \tau_1, \quad \tau_1, \tau_2 \in \mathbf{eps}\sigma\mathbb{Z} \quad \text{and} \quad |\tau_2| \leq \mathbf{eps} \cdot \text{ufp}(\tau_1),$$

$$(4.9) \quad |s + \varrho - \mathbf{res}| \leq 2\mathbf{eps}(1 - 2^{-M-1})\text{ufp}(\mathbf{res}).$$

If  $\sigma \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , then all components of the vector  $p'$  are zero and  $s + \varrho = \tau_1 + \tau_2$ .

If  $\mathbf{res} = 0$ , then  $s + \varrho = \tau_1 = \tau_2 = 0$  and all components of the vector  $p'$  are zero.

If  $\sigma > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , then

$$(4.10) \quad \text{ufp}(\tau_1) \geq 2^{2M+1}\mathbf{eps}\sigma.$$

REMARK 1. A standard application is the case  $\varrho = 0$  which will be our Algorithm 4.4 (**AccSum**). In this case, of course, also  $\varrho \in \mathbf{eps}\sigma_0\mathbb{Z}$ . We stated this piece of code with offset  $\varrho$  since we will use it for the calculation of a non-overlapping sequence of numbers representing  $K$ -fold accuracy (Algorithm 5.6).

REMARK 2. Besides the main result that  $\mathbf{res}$  is a faithful rounding, Lemma 4.3 also formulates some properties of the results of Algorithm 4.1 (**Transform**) independent of the (internal) final value of  $m$ .

PROOF OF LEMMA 4.3. For the final value of  $m$  in Algorithm 4.1 (**Transform**), (4.3), (4.2), (2.15) and (2.20) imply (4.8). Moreover,  $p' = p^{(m)}$ ,  $\sigma = \sigma_{m-1}$ , (4.1), (4.3) and (4.2) yield (4.7).

If  $\sigma \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  for the final value of  $m$  in Algorithm 4.1 (**Transform**), then (4.7) implies  $|p'_i| < \mathbf{eta}$ , so all components of the vector  $p'$  must be zero. Hence  $s + \varrho = \tau_1 + \tau_2$ , and  $\mathbf{res} = \text{fl}(\tau_1 + \tau_2) = \text{fl}(s + \varrho)$ , which is of course a faithful rounding, and (4.9) follows by (2.20).

Henceforth assume  $\sigma > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ . Then (4.4) implies (4.10), and we have to prove that  $\mathbf{res}$  is a faithful rounding of the exact result. Our assumptions imply

$$(4.11) \quad n + 2 \leq 2^M \quad \text{and} \quad 2^{2M+1}\mathbf{eps} \leq 1.$$

Next we abbreviate

$$(4.12) \quad \begin{aligned} \tau_3 &= \text{fl}(\sum_{i=1}^n p'_i) = \sum_{i=1}^n p'_i - \delta_3, \\ \tau'_2 &= \text{fl}(\tau_2 + \tau_3) = \tau_2 + \tau_3 - \delta_2, \\ \mathbf{res} &= \text{fl}(\tau_1 + \tau'_2) = \tau_1 + \tau'_2 - \delta_1. \end{aligned}$$

We will use Lemma 2.5 to prove that  $\mathbf{res}$  is a faithful rounding of  $s$ . By (4.7) and (4.12),  $s + \varrho = \tau_1 + \tau_2 + \tau_3 + \delta_3 = \tau_1 + \tau'_2 + \delta_2 + \delta_3$ , so

$$(4.13) \quad s + \varrho = r + \delta \quad \text{and} \quad \mathbf{res} = \text{fl}(r) \quad \text{for} \quad r := \tau_1 + \tau'_2 \quad \text{and} \quad \delta := \delta_2 + \delta_3.$$

By (4.7) and (2.16),

$$(4.14) \quad |\tau_3| = |\text{fl}(\sum_{i=1}^n p'_i)| \leq n\mathbf{eps}\sigma.$$

The standard estimation of floating-point sums [18] gives  $|\delta_3| \leq \gamma_{n-1} \sum_{i=1}^n |p'_i|$  independent of the order of summation, where  $\gamma_k := k\mathbf{eps}/(1 - k\mathbf{eps})$ . Note that this is also satisfied in the presence of underflow since summation is exact in that case. A little computation using  $n^2\mathbf{eps} < 2^{2M}\mathbf{eps} < 1$  shows  $\gamma_{n-1} < n\mathbf{eps}$ , and (4.7) gives

$$(4.15) \quad |\delta_3| < n^2\mathbf{eps}^2\sigma.$$

For later use we note that (4.14), (4.11) and (4.10) imply

$$(4.16) \quad |\tau_3| < 2^M\mathbf{eps}\sigma < 2^{-M-1}|\tau_1|.$$

Next (2.20), (4.12), (4.8) and (4.14) yield

$$(4.17) \quad |\delta_2| \leq \mathbf{eps}|\tau_2 + \tau_3| \leq \mathbf{eps}^2(\mathbf{ufp}(\tau_1) + n\sigma)$$

and (4.12), (4.8), (4.14) and (4.17) give

$$(4.18) \quad \begin{aligned} |\tau_1 + \tau'_2| &\geq |\tau_1 + \tau_2| - |\tau_2 - \tau'_2| = |\tau_1 + \tau_2| - |\tau_3 - \delta_2| \\ &\geq (1 - \mathbf{eps} - \mathbf{eps}^2)|\tau_1| - (1 + \mathbf{eps})n\mathbf{eps}\sigma, \end{aligned}$$

so (2.20) and (4.12) imply

$$(4.19) \quad |\mathbf{res}| \geq (1 - \mathbf{eps})|\tau_1 + \tau'_2| > (1 - 2\mathbf{eps})|\tau_1| - n\mathbf{eps}\sigma,$$

a lower bound for  $|\mathbf{res}|$ , also in the presence of underflow.

Now (4.13), (4.17) and (4.15) yield

$$(4.20) \quad |\delta| = |\delta_2 + \delta_3| < \mathbf{eps}^2(\mathbf{ufp}(\tau_1) + n\sigma + n^2\sigma).$$

Furthermore (4.20), (4.11), (4.10) and (4.19) imply

$$(4.21) \quad \begin{aligned} 2\mathbf{eps}^{-1}|\delta| &< 2\mathbf{eps} \cdot \mathbf{ufp}(\tau_1) + n(3 + 2n)\mathbf{eps}\sigma - n\mathbf{eps}\sigma \\ &< 2\mathbf{eps} \cdot \mathbf{ufp}(\tau_1) + (2^M - 2)2^{M+1}\mathbf{eps}\sigma - n\mathbf{eps}\sigma \\ &= 2\mathbf{eps} \cdot \mathbf{ufp}(\tau_1) + (1 - 2^{-M+1})2^{2M+1}\mathbf{eps}\sigma - n\mathbf{eps}\sigma \\ &\leq 2\mathbf{eps} \cdot \mathbf{ufp}(\tau_1) + (1 - 4\mathbf{eps} - 2^{-M})\mathbf{ufp}(\tau_1) - n\mathbf{eps}\sigma \\ &= (1 - 2\mathbf{eps} - 2^{-M})\mathbf{ufp}(\tau_1) - n\mathbf{eps}\sigma \\ &\leq (1 - 2\mathbf{eps})(1 - 2^{-M})\mathbf{ufp}(\tau_1) - (1 - 2^{-M})n\mathbf{eps}\sigma \\ &< (1 - 2^{-M})|\mathbf{res}| < |\mathbf{res}|. \end{aligned}$$

Using (4.13) and (2.20) we conclude

$$|s + \varrho - \mathbf{res}| \leq |\mathbf{fl}(r) - r| + |\delta| \leq \mathbf{eps} \cdot \mathbf{ufp}(\mathbf{res}) + \frac{1}{2}\mathbf{eps}(1 - 2^{-M})|\mathbf{res}|,$$

and  $|\mathbf{res}| < 2\mathbf{ufp}(\mathbf{res})$  proves (4.9).

If  $|\tau_1| \geq \mathbf{eps}^{-1}\mathbf{eta}$ , then (4.19), (4.11) and (4.10) yield

$$|\mathbf{res}| > (1 - 2\mathbf{eps} - 2^{-M-1})|\tau_1| > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}.$$

That means  $\mathbf{res} \notin \mathbb{U}$ , and (4.13), (4.21) and Lemma 2.5 show that  $\mathbf{res}$  is a faithful rounding of the sum  $s + \varrho$ . This leaves us with the case  $|\tau_1| < \mathbf{eps}^{-1}\mathbf{eta}$ . In that case (4.17), (4.8) and (4.16) yield

$$(4.22) \quad |\delta_2| \leq \mathbf{eps}|\tau_2 + \tau_3| \leq \mathbf{eps}^2|\tau_1| + 2^{-M-1}\mathbf{eps}|\tau_1| < \frac{1}{2}\mathbf{eps}|\tau_1| < \mathbf{eta},$$

and (4.15), (4.10) and (4.11) imply

$$(4.23) \quad |\delta_3| < n^2\mathbf{eps} \cdot 2^{-2M-1}|\tau_1| < \frac{1}{2}\mathbf{eps}|\tau_1| < \mathbf{eta}.$$

But  $\delta_2$  and  $\delta_3$  are the errors of the single floating-point additions to compute  $\tau'_2$  and  $\tau_3$ , respectively, and by (2.4) they are floating-point numbers, so by (4.22) and (4.23) both must be zero. Hence (4.13) implies that  $\mathbf{res} = \mathbf{fl}(r) = \mathbf{fl}(s + \varrho)$  is a faithful rounding of the exact result  $s + \varrho$ .

Finally, let  $\mathbf{res} = 0$ . To establish a contradiction assume  $\sigma > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ . The “until”-condition yields  $|\tau_1| = |t^{(m)}| \geq 2^{2M+1}\mathbf{eps}\sigma > 0$ , and (4.19) and (4.16) imply

$$|\mathbf{res}| > (1 - 2\mathbf{eps})|\tau_1| - n\mathbf{eps}\sigma > \frac{1}{2}|\tau_1| > 0,$$

a contradiction. Therefore  $\sigma \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  and we already proved  $p'_i = 0$  for all  $i$ . Hence  $\tau_1 = \mathbf{fl}(\tau_1 + \tau_2) = \mathbf{fl}(s + \varrho) = 0$  and  $\tau_2 = 0$  because  $\mathbf{res} = 0$  is a faithful rounding of  $s + \varrho$ . The lemma is proved.  $\square$

A special application of Algorithm 4.1 (**Transform**) with  $\varrho \neq 0$  may be the following. Suppose it is known that one component,  $p_1$  say, of a vector  $p$  is much larger in magnitude than the others. Then the call **Transform**( $p(2 : n), p_1$ ), in Matlab notation, may reduce the number of loops since then  $\sigma_0$  only depends on the smaller components  $p_2, \dots, p_n$  and not on  $p_1$ . We will use this in Section 5, where we will present Algorithm 5.6 (**AccSumK**) to compute a faithfully rounded result of  $K$ -fold accuracy.

We are now in the position to state our first algorithm for computing a faithfully rounded result of the sum of a vector of floating-point numbers.

ALGORITHM 4.4. *Accurate summation with faithful rounding.*

```
function res = AccSum(p)
    [\tau_1, \tau_2, p'] = Transform(p, 0)
    res = fl(\tau_1 + (\tau_2 + (\sum_{i=1}^n p'_i)))
```

PROPOSITION 4.5. *Let  $p$  be a vector of  $n$  floating-point numbers and assume  $2(n+2)^2\mathbf{eps} \leq 1$ . Let  $\mathbf{res}$  be the result of Algorithm 4.4 (**AccSum**) applied to  $p$ . Then  $\mathbf{res}$  is a faithful rounding of  $s := \sum_{i=1}^n p_i$ .*

REMARK. Algorithm 4.4 (**AccSum**) is identical to the piece of code we analyzed in Lemma 4.3, only the offset  $\varrho$  is set to zero and the output parameter  $\sigma$  in **Transform**, which is unnecessary here, is omitted.

PROOF. For zero input vector  $p$ , Algorithm 4.1 (**Transform**) implies  $\tau_1 = \tau_2 = p'_i = \mathbf{res} = 0$  for all  $i$ . For nonzero input vector  $p$ , the assumptions of Lemma 4.3 are satisfied, and the assertion follows.  $\square$

Lemma 2.6 implies the following corollary.

COROLLARY 4.6. *Under the assumption of Proposition 4.5 the computed result  $\mathbf{res}$  of Algorithm 4.4 (**AccSum**) is equal to the exact result  $s = \sum p_i$  if  $s$  is a floating-point number, or if  $\mathbf{res} \in \mathbb{U}$ . Moreover,*

$$|s| \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta} \quad \Rightarrow \quad \mathbf{res} \in \mathbb{U} \quad \Rightarrow \quad \mathbf{res} = s .$$

Therefore especially  $\mathbf{res} = 0$  if and only if  $s = 0$ , and

$$\mathbf{sign}(\mathbf{res}) = \mathbf{sign}(s) .$$

REMARK. The exact determination of the sign of a sum by Algorithm 4.4 is critical in the evaluation of geometrical predicates [19, 2, 43, 7, 27, 6, 12]. Rewriting a dot product as a sum by Algorithm 1.1 (**TwoProduct**), we can determine the exact sign of a dot product as well, which in turn decides whether a point is exactly on or on which side of a plane it is.

We showed that the result of Algorithm 4.4 is *always* a faithful rounding of the exact sum. Computational evidence suggests that the cases, where the result of **AccSum** is not rounded to nearest, are very rare. In several billion tests we never encountered such a case. This behavior can be explained as follows.

By (4.12) the error  $s + \varrho - \mathbf{res}$  computes to  $\delta_1 + \delta_2 + \delta_3$ , comprising of the rounding  $\delta_1$  of  $\mathbf{res}$  plus the additional error term  $\delta = \delta_2 + \delta_3$ . So the question is, when is  $\delta$  large enough to change the rounding of  $\mathbf{res}$  so that it is not to nearest (but still faithful). From (4.20) we know that  $|\delta|$  is not larger than  $n^2\mathbf{eps}^2|\tau_1|$  plus some very small term, and only for huge values of  $n$  we may expect  $\delta$  to have an impact on the rounding of  $\mathbf{res}$ .

For example, for  $n = 1000$  that means that the true value of  $\delta_1$  must be as close as  $10^6 \mathbf{eps} \sim 10^{-10}$  to the switching point for rounding. One may expect that this happens once in some 10 billion cases.

However, we can construct examples with faithful but not rounding to nearest. Consider  $p = [1 \ \mathbf{eps} \ \mathbf{eps}^2]$ . Then  $\text{AccSum}(p)$  produces  $\tau_1 = 1$ ,  $\tau_2 = 0$  and  $p' = [0 \ \mathbf{eps} \ \mathbf{eps}^2]$ , and  $\mathbf{res} = 1$ . This is because IEEE 754 rounds tie to even, so  $\text{fl}(1 + \mathbf{eps}) = 1$ .

Changing the strict into an “almost always” rounded to nearest offers quite a reward, namely the computational effort of Algorithm 4.4 ( $\text{AccSum}$ ) depends solely on the condition number: only the more difficult the problem, the more time must be spent. The maximum number of iterations  $m$  can be estimated as follows. The condition number of summation for  $\sum p_i \neq 0$  is defined [18] by

$$\text{cond} \left( \sum p_i \right) := \limsup_{\varepsilon \rightarrow 0} \left\{ \left| \frac{\sum \tilde{p}_i - \sum p_i}{\varepsilon \sum p_i} \right| : |\tilde{p}| \leq \varepsilon |p| \right\},$$

where absolute value and comparison of vectors is to be understood componentwise. Obviously

$$(4.24) \quad \text{cond} \left( \sum p_i \right) = \frac{\sum |p_i|}{|\sum p_i|}.$$

The following theorem estimates the maximal number  $m$  of iterations needed in  $\text{AccSum}$  depending on the number of elements  $n$  and the condition number.

**THEOREM 4.7.** *Assume Algorithm 4.4 ( $\text{AccSum}$ ) is applied to a vector of floating-point numbers  $p_i$ ,  $1 \leq i \leq n$ , with nonzero sum  $s$  and  $2(n+2)^2 \mathbf{eps} \leq 1$ . Then the following is true. If*

$$(4.25) \quad \text{cond} \left( \sum p_i \right) \leq 2^{-2M-2} (1 - 2^{-M}) [2^{-M} \mathbf{eps}^{-1}]^m,$$

then Algorithm 4.1 ( $\text{Transform}$ ) called by Algorithm 4.4 ( $\text{AccSum}$ ) stops after at most  $m$  executions of the “repeat-until”-loop. If the “repeat-until”-loop in Algorithm 4.1 ( $\text{Transform}$ ) has been executed at least  $m \geq 2$  times, then

$$(4.26) \quad \text{cond} \left( \sum p_i \right) > 2^{-2M-2} (1 - 2^{-M}) [2^{-M} \mathbf{eps}^{-1}]^{m-1}.$$

If the repeat-until loop is executed  $m$  times and absolute value and comparison is counted as one flop, then Algorithm 4.4 ( $\text{AccSum}$ ) needs  $(4m+3)n + \mathcal{O}(1)$  flops.

**PROOF.** For the called Algorithm 4.1 ( $\text{Transform}$ ) we have  $p_i^{(0)} := p_i$ , and its initialization implies

$$(4.27) \quad n+2 \leq 2^M, \quad \max_i |p_i| = \mu, \quad \frac{1}{2} \sigma_0 < 2^M \mu \leq \sigma_0 \quad \text{and} \quad \sigma_k = \varphi^k \sigma_0$$

for  $\varphi := 2^M \mathbf{eps}$  and  $0 \leq k \leq m$ .

To establish a contradiction assume Algorithm 4.1 is not finished after  $m$  executions of the “repeat-until”-loop. Then the “until”-condition, which is not satisfied for the value  $m$ , implies  $\sigma_{m-1} > \frac{1}{2} \mathbf{eps}^{-1} \mathbf{eta}$  so that  $\text{fl}(2^{2M+1} \mathbf{eps} \sigma_{m-1}) = 2^{2M+1} \mathbf{eps} \sigma_{m-1}$ , and

$$(4.28) \quad |t^{(m)}| < 2^{2M+1} \mathbf{eps} \sigma_{m-1} = 2^{M+1} \sigma_m.$$

Since the “repeat-until”-loop is to be executed again, we conclude

$$(4.29) \quad s = t^{(m)} + \sum p_i^{(m)}$$

as in (4.6) in the proof of Lemma 4.2. Denote the condition number by  $C$ . Then (4.24) and (4.27) yield  $|s| = C^{-1} \sum |p_i| \geq \mu C^{-1}$ . Combining this with (4.29), (4.2), (4.27), (4.25) and using  $(1 - 2^{-M})^{-1} > 1 + 2^{-M-1}$

TABLE 4.1  
Minimum treatable condition numbers by Algorithm 4.4

| $n$    | $m = 1$          | $m = 2$             | $m = 3$             | $m = 4$             | $m = 5$             |
|--------|------------------|---------------------|---------------------|---------------------|---------------------|
| 100    | $1.1 \cdot 10^9$ | $7.5 \cdot 10^{22}$ | $5.3 \cdot 10^{36}$ | $3.7 \cdot 10^{50}$ | $2.6 \cdot 10^{64}$ |
| 1000   | $2.1 \cdot 10^6$ | $1.8 \cdot 10^{19}$ | $1.6 \cdot 10^{32}$ | $1.4 \cdot 10^{45}$ | $1.2 \cdot 10^{58}$ |
| 10000  | $5.1 \cdot 10^2$ | $2.8 \cdot 10^{14}$ | $1.5 \cdot 10^{26}$ | $8.5 \cdot 10^{37}$ | $4.6 \cdot 10^{49}$ |
| $10^5$ |                  | $6.8 \cdot 10^{10}$ | $4.7 \cdot 10^{21}$ | $3.2 \cdot 10^{32}$ | $2.2 \cdot 10^{43}$ |
| $10^6$ |                  | $1.6 \cdot 10^7$    | $1.4 \cdot 10^{17}$ | $1.2 \cdot 10^{27}$ | $1.0 \cdot 10^{37}$ |
| $10^7$ |                  | $2.5 \cdot 10^2$    | $1.2 \cdot 10^{11}$ | $7.3 \cdot 10^{19}$ | $3.9 \cdot 10^{28}$ |

TABLE 4.2  
Minimum treatable length  $n$  for condition number  $\mathbf{eps}^{-1}$

| $m = 2$ | $m = 3$          | $m = 4$          |
|---------|------------------|------------------|
| 4096    | $1.0 \cdot 10^6$ | $6.7 \cdot 10^7$ |

implies

$$\begin{aligned}
 |t^{(m)}| &\geq |s| - \left| \sum p_i^{(m)} \right| \geq \mu C^{-1} - 2^M \mathbf{eps} \sigma_{m-1} \\
 &= \mu C^{-1} - \varphi^m \sigma_0 > (2^{-M-1} C^{-1} - \varphi^m) \sigma_0 \\
 &\geq (2^{-M-1} 2^{2M+2} (1 + 2^{-M-1}) - 1) \varphi^m \sigma_0 = 2^{M+1} \varphi^m \sigma_0 \\
 &= 2^{M+1} \sigma_m,
 \end{aligned}$$

a contradiction to (4.28). This proves our first assertion.

Next assume the “repeat-until”-loop in Algorithm 4.1 (**Transform**) has been executed at least  $m \geq 2$  times. Then the “until”-condition was not satisfied for  $m - 1$ , so that  $|t^{(m-1)}| < 2^{2M+1} \mathbf{eps} \sigma_{m-2} = 2^{M+1} \sigma_{m-1}$ . Hence (4.1), (4.2), (4.27) and  $1 + 2^{-M-1} + \frac{1}{2} \mathbf{eps} < (1 - 2^{-M})^{-1}$  imply

$$|s| = |t^{(m-1)} + \tau^{(m)} + \sum_{i=1}^n p_i^{(m)}| < (2^{M+1} + 1 + 2^M \mathbf{eps}) \sigma_{m-1} < 2^{M+1} (1 - 2^{-M})^{-1} \sigma_{m-1}.$$

Hence (4.24) and (4.27) yield

$$C \geq \frac{\mu}{|s|} > \frac{2^{-M-1} \sigma_0}{2^{M+1} (1 - 2^{-M})^{-1} \sigma_{m-1}} = 2^{-2M-2} (1 - 2^{-M}) \varphi^{-m+1}.$$

This proves the second assertion. Up to order 1, the calculation of  $\mu$  requires  $2n$  flops, **ExtractVector** requires  $4n$  and the computation of **res** requires  $n$  flops.  $\square$

For IEEE 754 double precision, Theorem 4.7 basically means that at least for condition numbers up to

$$\text{cond} \left( \sum p_i \right) \lesssim 2^{m(53-M)-2M-2}$$

Algorithm 4.4 (**AccSum**) computes a faithfully rounded result in at most  $m$  executions of the repeat-until loop. In Table 4.1 we show the lower bound by Theorem 4.7 for the condition number which can be treated for different values of  $n$  and  $m$ , where treatable means to produce a faithfully rounded result.

Conversely, we can use Theorem 4.7 to compute for given  $m$  and condition number  $\mathbf{eps}^{-1}$  the minimum length  $n$  of a vector for which a faithfully rounded result of its sum is computed. The value  $\mathbf{eps}^{-1}$  is the condition number for which we cannot expect a single correct digit by traditional recursive summation.

The value  $6.7 \cdot 10^7$  in Table 4.2 corresponds to the maximum value of  $n$  satisfying  $2(n+2)^2 \mathbf{eps} \leq 1$ . This means for condition number up to  $\mathbf{eps}^{-1}$  never more than  $m = 4$  iterations are needed. Note that the lower

TABLE 4.3  
*Floating-point operations needed for different dimension and condition number*

|                                    | AccSum | Sum2 | XBLAS |
|------------------------------------|--------|------|-------|
| $n = 1000, \text{ cond} = 10^6$    | $7n$   | $7n$ | $10n$ |
| $n = 1000, \text{ cond} = 10^{16}$ | $10n$  | $7n$ | $10n$ |
| $n = 10^6, \text{ cond} = 10^{16}$ | $13n$  | $7n$ | $10n$ |

and upper bound for the condition number in Theorem 4.7 differ by a factor  $2^M \mathbf{eps}$ , just the shift for  $\sigma$  in **Transform**. Algorithms for even larger values  $n > 6.7 \cdot 10^7$  will be presented in Section 7.

Our computing times compare favorably to other algorithms. Consider Algorithms 4.4 (**Sum2**) in [36] and the XBLAS summation algorithm **BLAS\_dsum\_x** in [31]. The results of both algorithms are of the same quality, namely as if computed in 2-fold precision. That means, for condition numbers up to  $\mathbf{eps}^{-1}$  we can expect a result accurate to the last bit. In Table 4.3 the required floating point operations are displayed for different vector lengths and condition numbers. Note that Algorithm 4.4 (**AccSum**) always computes a faithfully rounded result, independent of the condition number.

There is another big advantage of our Algorithm 4.4 (**AccSum**) compared to **Sum2** and XBLAS summation. Both **Sum2** and XBLAS summation pass through the input vector using the error-free transformation **TwoSum** (Algorithm 2.1), so that step  $i + 1$  of the transformation depends on step  $i$ . Hence parallelization of the algorithms is difficult. In contrast, the inner loop of **AccSum** can be perfectly executed in parallel (see also the Matlab code Algorithm 10.1 in the Appendix).

Next we take extra care about zero sums. The algorithm for provably computing a rounding to nearest result of the sum must be postponed to Section 6 because we need results from the following Section 5 on algorithms with  $K$ -fold accuracy.

We can significantly improve Algorithm 4.1 (**Transform**) by checking for  $t^{(m)} = 0$ . In that case  $\mathbf{fl}(t^{(m-1)} + \tau^{(m)}) = 0 = t^{(m-1)} + \tau^{(m)}$  because addition in the underflow range is exact, so (4.1) implies  $s + \varrho = \sum_{i=1}^n p_i^{(m)}$ . Hence, we can apply Algorithm 4.1 recursively to  $p_i^{(m)}$ .

This has two advantages. First, if  $s = 0$ , then  $p_i^{(m)}$  consists only of zeros at a certain stage. With checking for  $t^{(m)} = 0$  this is determined immediately in the recursive call of **Transform** by  $\mu = 0$ . Without this check, the “while”-loop would continue until  $\sigma_{m-1}$  is decreased into the underflow range, each time processing the whole vector  $p^{(m-1)}$  which may long be zero. Second, the input vector  $p_i^{(0)}$  might consist of some components large in absolute value which cancel exactly to zero, and some remaining small components. Then the latter form the new input vector  $p_i^{(m)}$ , and the recursive call skips the gap in one step.

The following modification of Algorithm 4.1 (**Transform**) includes this check for zero. If  $t^{(m)} = 0$  it may be advantageous to eliminate zero components in the vector  $p_i^{(m)}$ . This is done in the executable Matlab code given in the Appendix.

ALGORITHM 4.8. *Modified Algorithm 4.1 (Transform) with check for zero.*

```

function  $[\tau_1, \tau_2, p, \sigma, M] = \text{Transform}(p, \varrho)$ 
     $\mu = \max(|p_i|)$ 
    if  $\mu = 0$ ,  $\tau_1 = \varrho$ ,  $\tau_2 = 0$ , return, end if
     $M = \lceil \log_2(\text{length}(p^{(0)}) + 2) \rceil$ 
     $\sigma' = 2^{M + \log_2(\mu)}$ 
     $t' = \varrho$ 
    repeat
         $t = t'$ ;  $\sigma = \sigma'$ 
         $[\tau, p] = \text{ExtractVector}(\sigma, p)$ 
         $t' = \text{fl}(t + \tau)$ 
        if  $t' = 0$ ,  $[\tau_1, \tau_2, p] = \text{Transform}(p, 0)$ ; return; end if
         $\sigma' = \text{fl}(2^M \text{eps} \sigma)$ 
    until  $\sigma \leq \frac{1}{2} \text{eps}^{-1} \text{eta}$  or  $|t'| \geq \text{fl}(2^{2M+1} \text{eps} \sigma)$ 
     $[\tau_1, \tau_2] = \text{FastTwoSum}(t, \tau)$ 

```

Compared to Algorithm 4.1, the modified Algorithm 4.8 omits indices and adds  $M$  as output, the latter being used in Section 7 on huge  $n$ . Note that for clarity we still use the logarithm rather than Algorithm 3.5 (NextPowerTwo), but replace this in the Matlab code given in the Appendix.

Otherwise, both algorithms are identical except the check for  $t' = 0$ . Therefore all assertions on Algorithm 4.1 given in Lemma 4.2, 4.3, etc. are true for Algorithm 4.8 as well. Hence we may safely use the same name **Transform** for both algorithms. Note that the results might in fact be different because the shift  $\sigma$  may change, however, the *properties* such as the faithful rounding of the result etc. remain unchanged.

The algorithm can still be slightly improved. For  $\sigma$  just before entering the underflow range one call to **ExtractVector** may be saved depending on  $t'$ .

Next we prepare for the computation of a result with faithfully rounded  $K$ -fold accuracy, stored in a result vector of length  $K$ . The aim is to assure that the result vector is a non-overlapping sequence of numbers.

**5.  $K$ -fold accuracy.** Sometimes the precision of the given floating-point format may not be sufficient. One possibility for improvement is the use of multiple precision arithmetic such as [3, 5, 13, 1]. Frequently such packages support a special long precision data type.

An alternative may be to put together a long number by pieces of floating-point numbers, for example XBLAS uses this to simulate quadruple precision [31]. For more than two pieces this technique was used, for example, in [41] to compute highly accurate inclusions of the solution of systems of linear equations using a so-called long accumulator, which was developed in the 1970's in Karlsruhe for accurate computation of dot products. To gain as much accuracy as possible, the pieces should be non-overlapping, i.e. their bit representations should have no bit in common.

DEFINITION 5.1. *Two floating-point numbers  $f_1, f_2$  with  $|f_1| \geq |f_2|$  are called non-overlapping if*

$$(5.1) \quad \text{ufp}(f_2) < 2\text{eps} \cdot \text{ufp}(f_1) .$$

*A sequence of  $k \geq 1$  floating-point numbers  $f_1, \dots, f_k$  is called non-overlapping, if  $|f_1| \geq |f_2| \geq \dots \geq |f_{k'}| \geq 0 = f_{k'+1} = \dots = f_k$  and  $f_i, f_{i+1}$  is non-overlapping for  $1 \leq i < k'$ . This includes a sequence of zeros.*

Non-overlapping sequences have been used for example by Priest [38, 39] in the realm of geometrical predicates. Shewchuck [44] observed that enforcing non-overlapping may require substantial effort. Therefore he

defined and used weak non-overlapping, which means that bits of  $f_2$  may overlap zero bits “at the end” of  $f_1$ .

Though more efficiently computable than Priest’s, weakly non-overlapping sequences may comprise (and do) of  $f_i$  consisting only of one single bit, the leading one. Thus we cannot expect a  $K$ -fold accuracy for a sequence of length  $K$ . This changes for (strongly) non-overlapping sequences as demonstrated in the following lemma.

LEMMA 5.2. *Let  $f_1, \dots, f_k$  be a non-overlapping sequence of floating-point numbers. Then*

$$(5.2) \quad \text{ufp}(f_k) \leq \mathbf{eps}^{k-1} \cdot \text{ufp}(f_1) .$$

Moreover,

$$(5.3) \quad \text{ufp}\left(\sum_{i=1}^k f_i\right) \leq \text{ufp}(f_1) .$$

REMARK. Note that Lemma 5.2 is applicable to any subsequence  $f_{i_1}, \dots, f_{i_m}$  with  $i_1 > \dots > i_m$  because it is non-overlapping as well.

PROOF. Let  $f, g$  be a non-overlapping sequence. Then Definition 5.1 implies

$$(5.4) \quad \text{ufp}(g) \leq \mathbf{eps} \cdot \text{ufp}(f) ,$$

and (5.2) follows. To prove (5.3) we first show

$$(5.5) \quad |f| + 2\text{ufp}(g) \leq 2\text{ufp}(f)$$

If  $f \in \mathbb{U}$ , then (2.1) and (5.2) imply  $\text{ufp}(g) \leq \frac{1}{2}\mathbf{eta}$  and therefore  $g = 0$ . If  $f \notin \mathbb{U}$ , then  $|f| \leq \text{pred}(2\text{ufp}(f)) = 2(1 - \mathbf{eps})\text{ufp}(f)$ , and (5.4) implies (5.5). Without loss of generality we can omit remainder zeros in the sequence  $f_1, \dots, f_k$  and assume  $f_k \neq 0$ . Then

$$\left| \sum_{i=1}^k f_i \right| < \sum_{i=1}^{k-1} |f_i| + 2\text{ufp}(f_k) \leq \sum_{i=1}^{k-2} |f_i| + 2\text{ufp}(f_{k-1}) \leq \dots \leq 2\text{ufp}(f_1) ,$$

and (5.3) follows. □

We are aiming on non-overlapping sequences in the sense of Definition 5.1. Note that this is a little more than a faithfully rounded result in  $K$ -fold precision because there may be “gaps” between adjacent members of a sequence.

Suppose  $f_1 \in \mathbb{F}$  is a faithful rounding of  $r \in \mathbb{R}$ , and let  $f_2 \in \mathbb{F}$  be a faithful rounding of  $r - f_1$ . Only from the definition of faithful rounding this is not enough to ensure non-overlapping of  $f_1$  and  $f_2$ . Consider  $r = 1 - \mathbf{eps}^3$ . Then

$$f_1 := 1 - \mathbf{eps} = \text{pred}(1) \in \square(r) \quad \text{and} \quad f_2 := \mathbf{eps} \in \square(\mathbf{eps} - \mathbf{eps}^3) ,$$

but  $f_2$  is the last bit of  $f_1$ . Fortunately, our Algorithm 4.4 (AccSum) produces a faithfully rounded result good enough to guarantee non-overlapping.

LEMMA 5.3. *Let the assumptions of Lemma 4.3 be valid so that some  $\mathbf{res} \in \mathbb{F}$  has been computed with the code given over there and  $\mathbf{res} \in \square(s + \varrho)$ . Furthermore, suppose  $f_2 \in \square(s + \varrho - \mathbf{res})$  for some  $f_2 \in \mathbb{F}$ . Then  $\mathbf{res}, f_2$  are non-overlapping.*

PROOF. If  $\mathbf{res} \in \mathbb{U}$ , then Lemma 2.6 implies  $\mathbf{res} = s + \varrho$  and  $f_2 = 0$ . Therefore, we may assume  $\mathbf{res} \notin \mathbb{U}$ , so that  $f_2 \in \square(s + \varrho - \mathbf{res})$ , (5.1), (2.31) and (4.9) give

$$\begin{aligned} |f_2| &\leq |s + \varrho - \mathbf{res}| + |s + \varrho - \mathbf{res} - f_2| \leq |s + \varrho - \mathbf{res}| + 2\mathbf{eps} \cdot \text{ufp}(s + \varrho - \mathbf{res}) \\ &\leq (1 + 4\mathbf{eps})|s + \varrho - \mathbf{res}| \leq 2\mathbf{eps}(1 + 4\mathbf{eps})(1 - 2^{-M-1})\text{ufp}(\mathbf{res}) \\ &< 2\mathbf{eps} \cdot \text{ufp}(\mathbf{res}), \end{aligned}$$

proving by (5.1) that  $\mathbf{res}$  and  $f_2$  are non-overlapping.  $\square$

When applying the code given in Lemma 4.3 to a vector  $p$  and  $\varrho \in \mathbb{F}$ , it calculates a faithful rounding  $\mathbf{res}$  of  $s + \varrho$ , and also extracts  $p$  into some vector  $p'$ . But we have no *equation* relating  $s$ ,  $\varrho$ ,  $\mathbf{res}$  and  $p'$ , but only an *estimation* following by  $\mathbf{res} \in \square(s + \varrho)$ . The equation we know is  $s + \varrho = \tau_1 + \tau_2 + \sum p'_i$  as by (4.7).

However, in order to continue applying Lemma 5.3 to produce a non-overlapping sequence, we need a faithful rounding of  $s + \varrho - \mathbf{res}$ . Computing this by applying the code given in Lemma 4.3 to  $\varrho$  and the vector  $p$  appended by  $-\mathbf{res}$  would extract the entire vector  $p$  again and be highly ineffective. Fortunately, the following Algorithm 5.4 (**TransformK**) gives an effective way to compute a single floating-point number  $R$  relating  $s$ ,  $\varrho$ ,  $\mathbf{res}$  and  $p'$  by an equation. It paves the way to compute a non-overlapping sequence of  $K$  numbers  $f_1, \dots, f_K$  approximating  $s$ , thus establishing an approximation of  $K$ -fold accuracy.

ALGORITHM 5.4. *Error-free vector transformation including faithful rounding.*

$$\begin{aligned} \text{function } [\mathbf{res}, R, p'] &= \text{TransformK}(p, \varrho) \\ [\tau_1, \tau_2, p', \sigma] &= \text{Transform}(p, \varrho) \\ \mathbf{res} &= \text{fl}(\tau_1 + (\tau_2 + (\sum_{i=1}^n p'_i))) \\ R &= \text{fl}(\tau_2 - (\mathbf{res} - \tau_1)) \end{aligned}$$

LEMMA 5.5. *Let  $p$  be a nonzero vector of  $n$  floating-point numbers and let  $\varrho \in \mathbb{F}$ . Abbreviate  $s := \sum p_i$ . Assume that  $2(n+2)^2\mathbf{eps} \leq 1$ , and assume  $\varrho \in \mathbf{eps}\sigma_0\mathbb{Z}$  is satisfied for  $M := \lceil \log_2(n+2) \rceil$ ,  $\mu := \max_i |p_i|$  and  $\sigma_0 = 2^{M+\lceil \log_2 \mu \rceil}$ . Let  $\mathbf{res}$ ,  $R$  and  $p'$  be the results of Algorithm 5.4 (**TransformK**) applied to  $p$  and  $\varrho$ . Then  $\mathbf{res}$  is a faithful rounding of  $s + \varrho$  and*

$$(5.6) \quad s + \varrho - \mathbf{res} = R + \sum_{i=1}^n p'_i.$$

*If  $\mathbf{res} = 0$ , then  $\mathbf{res} = s + \varrho = 0$  and  $R$  and all  $p'_i$  are zero. If the vector  $p'$  is nonzero, then*

$$(5.7) \quad R \in \mathbf{eps}\sigma'\mathbb{Z}$$

*is satisfied for  $\mu' := \max_i |p'_i|$  and  $\sigma' := 2^{M+\lceil \log_2 \mu' \rceil}$ .*

REMARK 1. The main part of the proof of Lemma 5.5 will be to show that no rounding error can occur in the computation of  $R$ , that is  $R = \tau_2 - (\mathbf{res} - \tau_1)$ . This proof is quite involved and moved to the Appendix.

REMARK 2. Algorithm **TransformK** replaces the pair  $\tau_1, \tau_2$  without error into the pair  $\mathbf{res}, R$ . Therefore it is suitable for cascading, see Algorithm 5.6 (**AccSumK**).

REMARK 3. As is Matlab convention, an output parameter, in this case  $M$  for **TransformK**, has been omitted.

PROOF OF LEMMA 5.5. Algorithm 5.4 (**TransformK**) uses the same piece of code as the one in Lemma 4.3 for the computation of  $\tau_1, \tau_2, p', \sigma$  and  $\mathbf{res}$ , so the assertions of Lemma 4.3 are valid. If  $\mathbf{res} = 0$ , then  $s + \varrho = 0$  since  $\mathbf{res}$  is a faithful rounding of  $s + \varrho$ , and all  $p'_i$  are zero. Moreover,  $0 = s + \varrho = \tau_1 + \tau_2$  and  $\tau_1 = \text{fl}(\tau_1 + \tau_2) = 0 = \tau_2$ , so  $R = 0$  as well. It remains to prove (5.6) and (5.7).

First assume  $\sigma \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ . Then  $p'_i = 0$  for all  $i$  by Lemma 4.3, and (4.7) and (4.8) imply  $s + \varrho = \tau_1 + \tau_2$  and  $\mathbf{res} = \mathbf{fl}(\tau_1 + \tau_2) = \tau_1$ , so  $R = \tau_2$  and (5.6) follows. The case  $\sigma \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  is finished.

Henceforth we can assume  $\sigma > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ . We can use the analysis of Algorithm 4.1 (**Transform**) and the notation in the proof of Lemma 4.3, especially (4.12). Then by (2.20), (4.8) and (4.16) we have

$$(5.8) \quad |\tau'_2| = |\mathbf{fl}(\tau_2 + \tau_3)| \leq (1 + \mathbf{eps})|\tau_2 + \tau_3| < (1 + \mathbf{eps})(\mathbf{eps} + 2^{-M-1})|\tau_1| < |\tau_1| ,$$

so  $\mathbf{res} = \mathbf{fl}(\tau_1 + \tau'_2)$  and Lemma 2.8 imply

$$(5.9) \quad \Delta := \mathbf{res} - \tau_1 = \mathbf{fl}(\mathbf{res} - \tau_1) \in \mathbb{F} .$$

Next we have to show

$$(5.10) \quad R = \mathbf{fl}(\tau_2 - \Delta) = \tau_2 - \Delta$$

provided  $\sigma > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ . This proof is quite involved and left to the Appendix. Now (5.9) and (5.10) yield  $R = \tau_1 + \tau_2 - \mathbf{res}$ , and (4.7) gives

$$s + \varrho - \mathbf{res} = \tau_1 + \tau_2 + \sum_{i=1}^n p'_i - \mathbf{res} = R + \sum_{i=1}^n p'_i .$$

This proves (5.6). To see (5.7) note that (2.17), (4.10) and (2.11) imply  $\Delta = \mathbf{res} - \tau_1 \in \mathbf{eps} \cdot \mathbf{ufp}(\tau_1) \subseteq 2^{2M+1}\mathbf{eps}^2\sigma\mathbb{Z}$ . Furthermore, (4.8) and (4.11) yield  $\tau_2 \in \mathbf{eps}\sigma\mathbb{Z} \subseteq 2^{2M+1}\mathbf{eps}^2\sigma\mathbb{Z}$ , and (2.15) gives

$$R = \tau_2 - \Delta \in 2^{2M+1}\mathbf{eps}^2\sigma\mathbb{Z} .$$

But (4.7) and the definition of  $\sigma'$  imply  $\mathbf{eps}\sigma' \leq 2^M\mathbf{eps}^2\sigma$ , and (5.7) follows by (2.11).  $\square$

Next we can formulate an algorithm to compute a result of  $K$ -fold accuracy, stored in a non-overlapping result vector  $\mathbf{Res}$  of length  $K$ .

**ALGORITHM 5.6.** *Accurate summation with faithful rounding and result of  $K$ -fold accuracy.*

```

function  $\mathbf{Res} = \mathbf{AccSumK}(p, K)$ 
   $p^{(0)} = p, R_0 = 0$ 
  for  $k = 1 : K$ 
     $[\mathbf{Res}_k, R_k, p^{(k)}] = \mathbf{TransformK}(p^{(k-1)}, R_{k-1})$ 
    if  $\mathbf{Res}_k \in \mathbb{U}$ , return; end if
  end for

```

**PROPOSITION 5.7.** *Let  $p$  be a vector of  $n$  floating-point numbers, abbreviate  $s := \sum p_i$  and assume  $2(n+2)^2\mathbf{eps} \leq 1$ . Let  $1 \leq K \in \mathbb{N}$  be given, and let  $\mathbf{Res}$  be the vector of results of Algorithm 5.6 (**AccSumK**) applied to  $p$  and  $K$ . Denote the final value of  $k$  by  $K'$ . Then*

$$(5.11) \quad s = \sum_{\nu=1}^k \mathbf{Res}_\nu + R_k + \sum_{i=1}^n p_i^{(k)} \quad \text{for } 1 \leq k \leq K'$$

and

$$(5.12) \quad \mathbf{Res}_k \text{ is a faithful rounding of } s - \sum_{\nu=1}^{k-1} \mathbf{Res}_\nu \quad \text{for } 1 \leq k \leq K' ,$$

so that

$$(5.13) \quad \mathbf{pred}(\mathbf{Res}_k) < s - \sum_{\nu=1}^{k-1} \mathbf{Res}_\nu < \mathbf{succ}(\mathbf{Res}_k) \quad \text{for } 1 \leq k \leq K' .$$

Moreover,  $\mathbf{Res}_1, \dots, \mathbf{Res}_{K'}$  is a non-overlapping sequence of floating-point numbers, i.e. no bits in the binary representation of  $\mathbf{Res}_{k-1}$  and  $\mathbf{Res}_k$  overlap for  $k = 2, \dots, K'$ .

Abbreviate  $\hat{\mathbf{Res}} := \sum_{k=1}^{K'} \mathbf{Res}_k$ . Then  $s \in \mathbb{U}$  implies  $K' = 1$  and  $\hat{\mathbf{Res}} = s$ , and  $\mathbf{Res}_{K'} \in \mathbb{U}$  implies  $\hat{\mathbf{Res}} = s$ . If  $\mathbf{Res}_{K'} \notin \mathbb{U}$ , then  $K' = K$  and the following error estimates are satisfied:

$$(5.14) \quad |s - \hat{\mathbf{Res}}| < \frac{2}{1 - \mathbf{eps}} \mathbf{eps}^K |s|$$

and

$$(5.15) \quad |s - \hat{\mathbf{Res}}| < 2\mathbf{eps}^K \mathbf{ufp}(\mathbf{Res}_1) \leq 2\mathbf{eps}^K |\mathbf{Res}_1| .$$

REMARK 1. In fact, the sum  $\sum \mathbf{Res}_k$  is frequently of better accuracy than  $K$ -fold since zero bits in the binary representation of  $s$  “between”  $\mathbf{Res}_{k-1}$  and  $\mathbf{Res}_k$  are not stored.

REMARK 2. For  $k \geq 3$ , say, it may be advantageous to eliminate zero summands in the vectors  $p^{(k)}$ .

REMARK 3. Also note that the limited exponent range poses no problem to achieve  $K$ -fold accuracy. Since the exact result  $s$  is a sum of floating-point numbers, it follows  $s \in \mathbf{eta}\mathbb{Z}$ , and also  $s - \sum f_k \in \mathbf{eta}\mathbb{Z}$  for arbitrary floating-point numbers  $f_k$ . So for a certain  $K$  the exact sum  $s$  is stored in  $\sum \mathbf{Res}_k$ .

PROOF OF PROPOSITION 5.7. For zero input vector  $p$  the assertions are evident; henceforth we assume  $p$  to be nonzero.

For  $k = 0$  the assumptions of Lemma 5.5 are satisfied, so  $\mathbf{Res}_1$  is a faithful rounding of  $s = \sum p_i^{(0)}$ . Moreover, by (5.7) the assumptions of Lemma 5.5 are satisfied for  $k \geq 1$  as well, which proves (5.12) and (5.13). By construction and Lemma 5.3, the sequence  $\mathbf{Res}_1, \dots, \mathbf{Res}_K$  is non-overlapping.

Abbreviate  $s_k := \sum p_i^{(k)}$ . Then Lemma 5.5 and (5.6) imply

$$\begin{aligned} s &= s_0 &= \mathbf{Res}_1 + R_1 + s_1 \\ & &= \mathbf{Res}_1 + \mathbf{Res}_2 + R_2 + s_2 \\ & &= \dots \\ & &= \sum_{k=1}^{K'} \mathbf{Res}_k + R_{K'} + s_{K'} , \end{aligned}$$

which is (5.11). If  $s \in \mathbb{U}$ , then (5.12) and (2.30) imply  $\hat{\mathbf{Res}} = \mathbf{Res}_1 = s$ . If Algorithm 5.6 (`AccSumK`) ends with  $\mathbf{Res}_{K'} \in \mathbb{U}$ , then (2.30) yields  $\mathbf{Res}_{K'} = R_{K'-1} + s_{K'-1}$ , so  $s = \sum_{k=1}^{K'} \mathbf{Res}_k = \hat{\mathbf{Res}}$ .

Assume  $\mathbf{Res}_{K'} \notin \mathbb{U}$ . Then of course  $K' = K$ , and (5.12), (2.31) and (5.2) give

$$(5.16) \quad \mathbf{ufp}(s - \sum_{k=1}^K \mathbf{Res}_k) < 2\mathbf{eps} \cdot \mathbf{ufp}(\mathbf{Res}_k) \leq 2\mathbf{eps}^K \mathbf{ufp}(\mathbf{Res}_1) ,$$

proving (5.15). If  $\mathbf{ufp}(\mathbf{Res}_1) \leq \mathbf{ufp}(s)$ , then (5.14) follows as well. Otherwise,  $\mathbf{Res}_1 \in \square(s)$  implies  $\mathbf{pred}(\mathbf{Res}_1) < s < \mathbf{succ}(\mathbf{Res}_1)$ , so that  $\mathbf{ufp}(s) < \mathbf{ufp}(\mathbf{Res}_1)$  is only possible for  $|\mathbf{Res}_1| = \mathbf{ufp}(\mathbf{Res}_1)$ . In that case

$$(1 - \mathbf{eps})|\mathbf{Res}_1| = \mathbf{pred}(\mathbf{Res}_1) < s < |\mathbf{Res}_1| ,$$

so that  $\mathbf{ufp}(\mathbf{Res}_1) = |\mathbf{Res}_1| < (1 - \mathbf{eps})^{-1}|s|$ . The proof is finished.  $\square$

Now it becomes clear why so much effort was spent (see the Appendix) to prove (5.6) and (5.7): It is the key to cascade `TransformK` in Algorithm `AccSumK`. Especially for large gaps in the input vector,  $R_{k-1}$  in `AccSumK` may be large compared to  $p^{(k-1)}$ . In that case many case distinctions are necessary in the proof in the Appendix to assure that no rounding error can occur in the computation of  $R_k$  in `TransformK`.

**6. Rounding to nearest.** Now we have also the tools to design an algorithm computing the rounded to nearest result of the sum of a vector of floating-point numbers in the sense of IEEE 754.

ALGORITHM 6.1. *Accurate summation with rounding to nearest.*

```

function resN = NearSum(p)
    [res, R, p'] = TransformK(p, 0)
    [δ', R', p''] = TransformK(p', R)
    if δ' = 0, resN = res, return, end if
    res' = fl(res + sign(δ')eps|res|)
    if res' = res
        μ = sign(δ')eps|res|; res' = fl(res + 2sign(δ')eps|res|)
    else
        μ = fl((res' - res)/2)
    end if
    if |δ'| < |μ|, resN = res, return, end if
    if |δ'| > |μ|, resN = res', return, end if
    δ'' = TransformK(p'', R')
    if δ'' = 0, resN = fl(res + μ), return, end if
    if sign(δ'') = sign(μ), resN = res', return, end if
    resN = res

```

REMARK. Following Matlab convention the second and third output argument are omitted in the third call of `TransformK` because they are not needed.

The proof that `resN` is the rounded to nearest result of  $s$  is not difficult, though a little lengthy due to various case distinctions.

THEOREM 6.2. *Let  $p$  be a vector of  $n$  floating-point numbers and assume  $2(n+2)^2\text{eps} \leq 1$ . Let `resN` be the result of Algorithm 6.1 (`NearSum`). Then `resN` is the rounded to nearest exact sum  $s := \sum_{i=1}^n p_i$  in the sense of IEEE 754.*

PROOF. The assumptions of Lemma 5.5 are satisfied, so

$$(6.1) \quad s = \mathbf{res} + R + \sum_{i=1}^n p'_i = \mathbf{res} + \delta' + R' + \sum_{i=1}^n p''_i \quad \text{and} \quad \delta' \in \square(s - \mathbf{res}),$$

and Lemma 2.6 implies

$$(6.2) \quad \text{sign}(\delta') = \text{sign}(s - \mathbf{res}).$$

We will determine the “switching point”  $M$  for rounding to nearest determined by `res` and the sign of  $\delta'$ . Then we have to check whether the true sum  $s$  is less than or greater than or equal to  $M$ .

If  $\delta' = 0$ , the final value of `resN` is the correct value `res` =  $s$ . Moreover, if `res`  $\in \mathbb{U}$ , then Lemma 2.6 implies `res` =  $s$  and  $\delta' = 0$ , so that `resN` = `res` =  $s$ .

Henceforth we may assume  $\delta' \neq 0$  and `res`  $\notin \mathbb{U}$ . Then (6.2) yields

$$(6.3) \quad \begin{array}{ll} \delta' < 0 & \Rightarrow \quad N(\mathbf{res}) := \text{pred}(\mathbf{res}) < s < \mathbf{res} \\ \delta' > 0 & \Rightarrow \quad \mathbf{res} < s < \text{succ}(\mathbf{res}) =: N(\mathbf{res}). \end{array}$$

That means the “switching point” deciding the rounding to nearest of  $s$  is  $M := \frac{1}{2}(\mathbf{res} + N(\mathbf{res}))$ , with the neighbor  $N(\mathbf{res})$  defined by (6.3) depending on the sign of  $\delta'$ . We claim that

$$(6.4) \quad M = \mathbf{res} + \mu \quad \text{and} \quad \mathbf{res}' = N(\mathbf{res}).$$

To prove (6.4) first assume  $|\mathbf{res}| \neq \text{ufp}(\mathbf{res})$ . Then  $|N(\mathbf{res}) - \mathbf{res}| = 2\mathbf{eps} \cdot \text{ufp}(\mathbf{res})$  since  $\mathbf{res} \notin \mathbb{U}$ , and

$$|M - \mathbf{res}| = \left| \frac{1}{2}(N(\mathbf{res}) - \mathbf{res}) \right| = \mathbf{eps} \cdot \text{ufp}(\mathbf{res}) < \mathbf{eps}|\mathbf{res}| < 2\mathbf{eps} \cdot \text{ufp}(\mathbf{res}) = |N(\mathbf{res}) - \mathbf{res}| ,$$

i.e.  $\mathbf{res} + \text{sign}(\delta')\mathbf{eps}|\mathbf{res}|$  is in the open interval with endpoints  $M$ , which is the midpoint between  $\mathbf{res}$  and  $N(\mathbf{res})$ , and  $N(\mathbf{res})$ . Hence  $\mathbf{res}' = N(\mathbf{res}) \neq \mathbf{res}$ . Furthermore,  $\mathbf{res} \notin \mathbb{U}$  implies  $(\mathbf{res}' - \mathbf{res})/2 = (N(\mathbf{res}) - \mathbf{res})/2 \in \mathbb{F}$ , so that  $\mu = M - \mathbf{res}$ .

Next assume  $|\mathbf{res}| = \text{ufp}(\mathbf{res})$ , so that  $|\mathbf{res}| \geq \mathbf{eps}^{-1}\mathbf{eta}$  since  $\mathbf{res} \notin \mathbb{U}$ , and

$$\begin{aligned} |N(\mathbf{res}) - \mathbf{res}| &= 2\text{ufp}(\mathbf{res}) & \text{for } |N(\mathbf{res})| &= \text{succ}(|\mathbf{res}|) , \\ |N(\mathbf{res}) - \mathbf{res}| &= \text{ufp}(\mathbf{res}) & \text{for } |N(\mathbf{res})| &= \text{pred}(|\mathbf{res}|) \end{aligned}$$

In the first case,  $\mathbf{res} + \text{sign}(\delta')\mathbf{eps}|\mathbf{res}|$  is equal to the midpoint  $M$  between  $\mathbf{res}$  and  $N(\mathbf{res})$ , so rounding of tie to even implies  $\text{fl}(\mathbf{res} + \text{sign}(\delta')\mathbf{eps}|\mathbf{res}|) = \mathbf{res}$ . Therefore  $\mu = M - \mathbf{res}$  and  $\mathbf{res}'$  is corrected to  $N(\mathbf{res})$ .

In the second case,  $\mathbf{res} + \text{sign}(\delta')\mathbf{eps}|\mathbf{res}|$  is equal to  $N(\mathbf{res})$ , the floating-point number adjacent to  $\mathbf{res}$  towards zero. Therefore  $\mathbf{res}' = N(\mathbf{res})$ . If  $|\mathbf{res}| \geq 2\mathbf{eps}^{-1}\mathbf{eta}$ , then  $|\mathbf{res}' - \mathbf{res}| \geq 2\mathbf{eta}$  so that  $(\mathbf{res}' - \mathbf{res})/2 \in \mathbb{F}$ ,  $\mu = (\mathbf{res}' - \mathbf{res})/2$  and (6.4). But  $|\mathbf{res}| = \mathbf{eps}^{-1}\mathbf{eta}$  is not possible because in this case  $s$  cannot be strictly between  $N(\mathbf{res})$  and  $\mathbf{res}$  since  $|N(\mathbf{res}) - \mathbf{res}| = \mathbf{eta}$  and, as a sum of floating-point numbers,  $s \in \mathbf{eta}\mathbb{Z}$ . This proves (6.4).

Of course,  $\text{sign}(\mu) = \text{sign}(\delta')$ . It follows

$$(6.5) \quad \begin{aligned} |s - \mathbf{res}| < |\mu| &\Rightarrow \text{fl}(s) = \mathbf{res} && \text{and} \\ |s - \mathbf{res}| > |\mu| &\Rightarrow \text{fl}(s) = \mathbf{res}' && \text{and} \\ |s - \mathbf{res}| = |\mu| &\Rightarrow \text{fl}(s) = \text{fl}(\mathbf{res} + \mu) . \end{aligned}$$

By Lemma 5.5,  $\delta'$  is a faithful rounding of  $s - \mathbf{res}$ , so  $\mu \in \mathbb{F}$  and (2.29) imply that  $|\delta'| \leq |\mu| \Rightarrow |s - \mathbf{res}| \leq |\mu|$ , and these cases are handled correctly in Algorithm `NearSum`. In case  $|\delta'| = |\mu|$ , which means  $\delta' = \mu$ , we need more work to decide the sign of  $s - \mathbf{res} - \mu$ . By Lemma 2.6 and (6.1),  $\delta''$  is a faithful rounding of  $s - \mathbf{res} - \mu = R' + \sum_{i=1}^n p_i''$ . So  $\delta' = \mu$  and Lemma 2.6 imply

$$\begin{aligned} \text{sign}(\mu) &= \text{sign}(s - \mathbf{res}) && \text{and} \\ \text{sign}(\delta'') &= \text{sign}(s - \mathbf{res} - \mu) . \end{aligned}$$

Hence Algorithm `NearSum` handles the cases correctly according to (6.5).  $\square$

In contrast to Algorithm 4.4 (`AccSum`), the computing time of Algorithm 6.1 (`NearSum`) depends on the maximum exponent  $E$  of the vector entries: The maximum number of iterations in `Transform` can only be bounded by some  $m \lesssim (E - \log_2(\mathbf{eps}^{-1}\mathbf{eta})) / (\log_2(\mathbf{eps}^{-1} - M))$ , which is about  $m \lesssim (E + 1024) / (53 - M)$  for IEEE 754 double precision. This is independent of the condition number of the sum. However, it seems unlikely that this maximum is achieved in other than constructed examples. It can only be achieved if there are vector elements  $|p_i|$  near the overflow range, the binary representations of the other vector elements completely cover the whole exponent range from 1023 downto  $-1024$  without gap, and the exact sum cancels completely such that the exact result is exactly the midpoint between two adjacent floating-point numbers.

**7. Vectors of huge length.** In IEEE 754 double precision our summation algorithm computes a faithful rounding of the exact result for up to about  $6.7 \cdot 10^7$  summands. This should suffice for most practical purposes. In IEEE 754 single precision with  $\mathbf{eps} = 2^{-24}$ , the number of summands is restricted to  $n = 2894$ , which may be an obstacle for the application of our algorithms. This can be improved by rewriting of Algorithm 4.8 (`Transform`).

We first observe that the extraction in **Transform** works for large  $n$  as well if we only change the “until”-condition on  $|t^{(m)}|$ .

LEMMA 7.1. *Assume the “until”-condition in Algorithm 4.8 (**Transform**) is changed into*

$$(7.1) \quad \text{until } \sigma \leq \frac{1}{2}\text{eps}^{-1}\text{eta} \text{ or } |t'| \geq \sigma$$

and call that algorithm **TransformP**. Let  $\tau_1, \tau_2, p', \sigma$  be the results of Algorithm **TransformP** applied to a nonzero vector of floating-point numbers  $p_i, 1 \leq i \leq n$ , and  $\varrho \in \mathbb{F}$ . Assume that  $2^M \text{eps} \leq \frac{1}{2}$  and  $\varrho \in \text{eps}\sigma_0\mathbb{Z}$  is satisfied for  $M := \lceil \log_2(n+2) \rceil$ ,  $\mu := \max_i |p_i|$  and  $\sigma_0 = 2^{M+\lceil \log_2 \mu \rceil}$ . Denote  $s := \sum_{i=1}^n p_i$ .

Then Algorithm **TransformP** will stop, and

$$(7.2) \quad s + \varrho = \tau_1 + \tau_2 + \sum_{i=1}^n p'_i,$$

$$(7.3) \quad \max |p'_i| \leq \text{eps}\sigma, \quad \tau_1 = \text{fl}(\tau_1 + \tau_2) \quad \text{and} \quad \tau_1, \tau_2 \in \text{eps}\sigma\mathbb{Z}.$$

If  $\sigma \leq \frac{1}{2}\text{eps}^{-1}\text{eta}$ , then the vector  $p'$  is identically zero. If  $\sigma > \frac{1}{2}\text{eps}^{-1}\text{eta}$ , then

$$(7.4) \quad |\tau_1| \geq \sigma.$$

REMARK. Note that the computation of  $\sigma'$  may be afflicted with a rounding error if  $\sigma$  is in the underflow range  $\mathbb{U}$ . However, in this case the algorithm stops as in the original version and  $\sigma'$  is not used any more.

PROOF OF LEMMA 7.1. Carefully going through the proof of correctness of Algorithm 4.1 (**Transform**) in Lemma 4.2, we see that the assumptions of Theorem 3.4 are satisfied, and that the only places in the proof of (4.1), (4.2) and (4.3), where the stronger assumption  $2(n+2)^2 \text{eps} \leq 1$  is used, are firstly  $2^M \text{eps} < 1$  so that the algorithm stops, and secondly the deduction of  $|t^{(m-1)}| < \sigma_{m-2}$  in (4.5). However, the latter was in the induction argument, assuming that the corresponding “until”-condition was not satisfied. The changed “until”-condition (7.1) implies  $|t^{(m-1)}| < \sigma_{m-2}$  directly.

Algorithm 4.8 (**Transform**) is identical to Algorithm 4.1 (**Transform**) despite check for zero sum and omitting indices, so as before the assertions are not affected. The remaining assertions are merely a transcription.  $\square$

So far, so good. However, by (7.2) and (7.3), the error of  $\tau_1$  to the true result  $s + \varrho$  may be as large as  $\sum p'_i$ , which is bounded by  $n \text{eps}\sigma$ . For huge  $n$  this may come close to  $\sigma$ , the lower bound for  $\tau_1$ . Nevertheless, it suffices to determine the sign of  $s$ .

THEOREM 7.2. *Assume the “until”-condition in Algorithm 4.8 (**Transform**) is changed into (7.1) and call that algorithm **TransformP**. Let  $\tau_1, \tau_2, p', \sigma$  be the results of Algorithm **TransformP** applied to a nonzero vector of floating-point numbers  $p_i, 1 \leq i \leq n$ , and  $\varrho \in \mathbb{F}$ . Assume that  $2^M \text{eps} \leq \frac{1}{2}$  and  $\varrho \in \text{eps}\sigma_0\mathbb{Z}$  is satisfied for  $M := \lceil \log_2(n+2) \rceil$ ,  $\mu := \max_i |p_i|$  and  $\sigma_0 = 2^{M+\lceil \log_2 \mu \rceil}$ . Denote  $s := \sum_{i=1}^n p_i$ .*

Then Algorithm **TransformP** will stop, and

$$(7.5) \quad \text{sign}(s + \varrho) = \text{sign}(\text{res}).$$

PROOF. The assumptions of Lemma 7.1 are satisfied. Hence  $\sigma \leq \frac{1}{2}\text{eps}^{-1}\text{eta}$  implies  $\tau_1 = \text{fl}(\tau_1 + \tau_2) = \text{fl}(s + \varrho)$  and  $\text{sign}(s + \varrho) = \text{sign}(\tau_1) = \text{sign}(\text{res})$ . If  $\sigma > \frac{1}{2}\text{eps}^{-1}\text{eta}$ , then  $|\tau_1| \geq \sigma \geq \text{eps}^{-1}\text{eta}$ . Hence (7.3) gives

$$|\tau_2| + \left| \sum_{i=1}^n p'_i \right| \leq \text{eps}|\tau_1| + 2^M \text{eps}\sigma \leq \left( \text{eps} + \frac{1}{2} \right) |\tau_1|.$$

and (7.2) finishes the proof.  $\square$

This simplified version of Algorithm 4.4 (`AccSum`) works in single precision for dimensions up to  $8.3 \cdot 10^6$ , and in double precision for dimensions up to  $4.5 \cdot 10^{15}$ .

Next we improve the approximation by continuing to extract  $p'_i$  until the error term is small enough compared to  $\tau_1$ . This is done in the following Algorithm 7.3 (`AccSumHugeN`) to achieve faithful rounding.

ALGORITHM 7.3. *Accurate summation with faithful rounding for huge  $n$ .*

```

function res = AccSumHugeN(p)
    [ $\tau_1, \tau_2, q^{(0)}, \sigma, M$ ] = TransformP(p, 0)
    if  $\sigma \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , res =  $\tau_1$ , return, end if
     $k = 0$ ;  $\sigma_0 = \mathbf{fl}(2^M \mathbf{eps} \sigma)$ 
    repeat
         $k = k + 1$ 
        [ $\tau^{(k)}, q^{(k)}$ ] = ExtractVector( $\sigma_{k-1}, q^{(k-1)}$ )
         $\sigma_k = \mathbf{fl}(2^M \mathbf{eps} \sigma_{k-1})$ 
    until  $\sigma_{k-1} \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$  or  $\sigma_{k-1} < \mathbf{fl}(4\mathbf{eps}|\tau_1|)$ 
     $\tilde{\tau}^{(1)} = \mathbf{fl}(\tau^{(1)} + (\tau^{(2)} + \dots + (\tau^{(k)} + (\tau_2 + (\sum_{i=1}^n q_i^{(k)}))) \dots))$ 
    res =  $\mathbf{fl}(\tau_1 + \tilde{\tau}^{(1)})$ 
    
```

PROPOSITION 7.4. *Assume the “until”-condition in Algorithm 4.8 (`Transform`) is changed into (7.1) and call that algorithm `TransformP`. Let  $p$  be a vector of  $n$  floating-point numbers, assume  $8(n+2)\mathbf{eps} \leq 1$  and  $\mathbf{eps} \leq 1/128$ . Let **res** be the result of Algorithm 7.3 (`AccSumHugeN`) applied to  $p$ . Then **res** is a faithful rounding of  $s := \sum_{i=1}^n p_i$ .*

REMARK 1. For IEEE 754 single precision with  $\mathbf{eps} = 2^{-24}$  this limits the vector length to a little over 2 million rather than  $n = 2894$  for Algorithm `AccSum`. For IEEE 754 double precision the vector length is now limited to about  $1.1 \cdot 10^{15}$ .

REMARK 2. Note that the code given in the Appendix uses Algorithm 3.5 (`NextPowerTwo`) to compute  $2^M$  rather than  $M$ , so that `AccSumHugeN` is to be adapted accordingly.

PROOF OF PROPOSITION 7.4. Algorithm `TransformP` is called with  $\varrho = 0$ . So  $2^M \mathbf{eps} \leq \frac{1}{2}$  and Lemma 7.1 imply

$$(7.6) \quad s = \tau_1 + \tau_2 + \sum_{i=1}^n q_i^{(0)},$$

$$(7.7) \quad \max |q_i^{(0)}| \leq \mathbf{eps} \sigma,$$

$$(7.8) \quad \tau_1 = \mathbf{fl}(\tau_1 + \tau_2) \quad \text{and} \quad |\tau_2| \leq \mathbf{eps} |\tau_1|.$$

If  $\sigma$  computed by Algorithm `TransformP` is in the underflow range, then  $q^{(0)}$  the zero vector by Lemma 7.1, and  $\mathbf{fl}(s) = \mathbf{res}$  by (7.6).

Henceforth we may assume  $\sigma > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , so that (7.4) implies

$$(7.9) \quad |\tau_1| \geq \sigma \geq \mathbf{eps}^{-1}\mathbf{eta}.$$

Therefore no rounding error occurs in the second part of the “until”-condition, i.e.  $\mathbf{fl}(4\mathbf{eps}|\tau_1|) = 4\mathbf{eps}|\tau_1|$ .

By definition of  $\sigma_0$  and (7.7),  $\max_i |q_i^{(0)}| \leq 2^{-M} \sigma_0$ , so the assumptions of Theorem 3.4 are satisfied for the first call of `ExtractVector` for  $k = 1$ . Therefore,

$$s = \tau_1 + \tau_2 + \tau^{(1)} + \sum_{i=1}^n q_i^{(1)}, \quad |\tau^{(1)}| < \sigma_0 \quad \text{and} \quad \max_i |q_i^{(1)}| \leq \mathbf{eps} \sigma_0 = 2^{-M} \sigma_1.$$

If  $\sigma_0 \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , then the loop finishes and a possible rounding error in the computation of  $\sigma_1$  does no harm since  $\sigma_1$  is not used. If  $\sigma_0 > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , then  $\sigma_1$  is computed without rounding error and can safely be used.

Again, the assumptions of Theorem 3.4 are satisfied for the next call of `ExtractVector`, and repeating this argument shows

$$(7.10) \quad s = \tau_1 + \tau_2 + \sum_{\nu=1}^k \tau^{(\nu)} + \sum_{i=1}^n q_i^{(k)}, \quad |\tau^{(k)}| < \sigma_{k-1} \quad \text{and} \quad \max_i |q_i^{(k)}| \leq 2^{-M} \sigma_k$$

for  $k$  between 1 and its final value  $K$ . The same argument as before applies to possible rounding errors in the computation of  $\sigma_K$ . Until further notice assume  $\sigma_{K-1} > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , so that the “until”-condition yields

$$(7.11) \quad \sigma_{K-1} < 4\mathbf{eps}|\tau_1|.$$

Denote

$$(7.12) \quad \tau^{(K+1)} := \tau_2 + \left( \sum_{i=1}^n q_i^{(K)} \right) \quad \text{and} \quad \tilde{\tau}^{(K+1)} := \mathfrak{fl}(\tau^{(K+1)}) = \tau^{(K+1)} - \delta_{K+1},$$

$$(7.13) \quad \tilde{\tau}^{(k)} := \mathfrak{fl}(\tau^{(k)} + \tilde{\tau}^{(k+1)}) = \tau^{(k)} + \tilde{\tau}^{(k+1)} - \delta_k \quad \text{for } 1 \leq k \leq K.$$

This implies  $\tilde{\tau}^{(1)} = \tau^{(1)} + \tilde{\tau}^{(2)} - \delta_1 = \tau^{(1)} + \tau^{(2)} + \tilde{\tau}^{(3)} - \delta_2 - \delta_1 = \dots$ , so that

$$(7.14) \quad \tilde{\tau}^{(1)} = \sum_{k=1}^{K+1} \tau^{(k)} - \sum_{k=1}^{K+1} \delta_k = \sum_{k=1}^K \tau^{(k)} + \tau_2 + \left( \sum_{i=1}^n q_i^{(K)} \right) - \sum_{k=1}^{K+1} \delta_k.$$

Abbreviate  $\varphi := 2^M \mathbf{eps}$ . The assumptions imply

$$(7.15) \quad \varphi \leq \frac{1}{8}, \quad \mathbf{eps} \leq \frac{1}{128}, \quad M \geq 2 \quad \text{and} \quad \sigma_k = \varphi^{k+1} \sigma \quad \text{for } 0 \leq k \leq K.$$

Then (7.10), (7.8), the standard estimation for floating-point sums [18, Lemma 8.4] and (7.11) yield

$$(7.16) \quad |\delta_{K+1}| \leq \gamma_n (|\tau_2| + \sum_{i=1}^n |q_i^{(K)}|) \leq \frac{1}{7} (|\tau_2| + \sigma_K) \leq \frac{3}{14} \mathbf{eps} |\tau_1|.$$

Next we need an upper bound on  $|\tilde{\tau}^{(k)}|$ . We use an induction argument to show

$$(7.17) \quad |\tilde{\tau}^{(k)}| \leq \left( \sum_{\nu=0}^{K+1-k} (1 + \mathbf{eps})^{\nu+1} \varphi^\nu \right) \sigma_{k-1} + (1 + \mathbf{eps})^{K+2-k} |\tau_2| \quad \text{for } 1 \leq k \leq K+1.$$

For  $k = K+1$  this follows by

$$\begin{aligned} |\tilde{\tau}^{(K+1)}| &\leq |\mathfrak{fl}(\tau_2 + \left( \sum_{i=1}^n q_i^{(K)} \right))| \leq (1 + \mathbf{eps}) (|\tau_2| + |\mathfrak{fl}(\sum_{i=1}^n q_i^{(K)})|) \leq (1 + \mathbf{eps}) (|\tau_2| + n 2^{-M} \sigma_K) \\ &\leq (1 + \mathbf{eps}) (|\tau_2| + \sigma_K) \end{aligned}$$

using (7.12), (2.16) and (7.10), and the induction step uses

$$|\tilde{\tau}^{(k)}| \leq (1 + \mathbf{eps}) (|\tau^{(k)}| + |\tilde{\tau}^{(k+1)}|)$$

by (7.13), and  $|\tau^{(k)}| < \sigma_{k-1}$  by (7.10). Hence (7.15) gives

$$|\tilde{\tau}^{(k)}| \leq \frac{1 + \mathbf{eps}}{1 - (1 + \mathbf{eps})\varphi} \varphi^k \sigma + (1 + \mathbf{eps})^{K+2-k} |\tau_2| \quad \text{for } 1 \leq k \leq K+1,$$

and by (7.13) and (7.15) it follows

$$(7.18) \quad \mathbf{eps}^{-1} |\delta_k| \leq |\tilde{\tau}^{(k)}| \leq \frac{7}{6} \varphi^k \sigma + (1 + \mathbf{eps})^{K+2-k} |\tau_2| \quad \text{for } 1 \leq k \leq K.$$

Next we need an upper bound on  $K$ . Denote  $\mathbf{eps} = 2^{-m}$ , then (7.15) gives  $m \geq 7$  and  $M - m \leq -3$ . For  $K > 2$ , the “until”-condition and (7.9) imply  $\varphi^{K-1}\sigma = \sigma_{K-2} \geq 4\mathbf{eps}|\tau_1| \geq 4\mathbf{eps}\sigma$ , so that  $(K-1)(M-m) \geq 2-m$  and

$$(7.19) \quad K \leq 1 + \frac{m-2}{m-M} \leq \frac{m+1}{3} .$$

This is of course also true for  $K \leq 2$ . For later use we use  $m \geq 7$  to show

$$(7.20) \quad (1 + \mathbf{eps})^2 \frac{(1 + \mathbf{eps})^K - 1}{\mathbf{eps}} \leq \left(\frac{129}{128}\right)^2 ((1 + 2^{-m})^{\frac{m+1}{3}} - 1) \mathbf{eps}^{-1} < \frac{1}{28} \mathbf{eps}^{-1} .$$

Moreover,  $m \geq 7$  implies  $(1 + \mathbf{eps})^{K+1} \leq (1 + 2^{-m})^{1 + \frac{m+1}{3}} \leq 4/3$ , so that the definition of  $\mathbf{res}$ , (7.18), (7.8), (7.15) and (7.9) yield

$$(7.21) \quad \begin{aligned} |\mathbf{res}| &\geq (1 - \mathbf{eps})(|\tau_1| - |\tilde{\tau}^{(1)}|) \\ &\geq (1 - \mathbf{eps})(|\tau_1| - \frac{7}{6}\varphi\sigma - (1 + \mathbf{eps})^{K+1}\mathbf{eps}|\tau_1|) \\ &\geq (1 - \mathbf{eps})\left(\frac{95}{96}|\tau_1| - \frac{7}{48}\sigma\right) \\ &> \frac{1}{2}\sigma \geq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta} , \end{aligned}$$

proving that  $\mathbf{res} \notin \mathbb{U}$ . Moreover, (7.15) and (7.21) yield

$$(7.22) \quad |\mathbf{res}| > \frac{5}{6}|\tau_1| .$$

Now set  $r := \tau_1 + \tilde{\tau}^{(1)}$ , so that  $\mathbf{res} = \mathbf{fl}(r)$ . Then (7.10) and (7.14) imply

$$\delta := s - r = \sum_{k=1}^{K+1} \delta_k .$$

We will show  $2|\delta| < \mathbf{eps}|\mathbf{res}|$  and apply Lemma 2.5 to demonstrate that  $\mathbf{res}$  is a faithful rounding of the true result  $s$ . With (7.18) and using (7.15) it holds

$$(7.23) \quad \mathbf{eps}^{-1}|\delta| \leq \frac{7}{6} \frac{\varphi}{1-\varphi} \sigma + (1 + \mathbf{eps})^2 \frac{(1 + \mathbf{eps})^K - 1}{\mathbf{eps}} |\tau_2| + \mathbf{eps}^{-1} |\delta_{K+1}| .$$

Using (7.15), (7.8), (7.16), (7.20),  $\sigma \leq |\tau_1|$  and (7.22), a little computation shows

$$(7.24) \quad 2\mathbf{eps}^{-1}|\delta| < 2 \left( \frac{1}{6}\sigma + \frac{1}{4}|\tau_1| \right) \leq \frac{5}{6}|\tau_1| < |\mathbf{res}| .$$

By (7.21) we know  $\mathbf{res} \notin \mathbb{U}$ , so Lemma 2.5 proves that  $\mathbf{res}$  is a faithful rounding of  $s$ . This closes the case  $\sigma_{K-1} > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ .

Suppose  $\sigma_{K-1} \leq \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ , then (7.10) implies  $q_i^{(K)} = 0$  for all  $i$  and the second part of the “until”-condition, that is (7.11) need not be true. But both are only needed to derive the upper bound of  $\delta_{K+1}$ , and (7.12) implies  $\delta_{K+1} = 0$ . Especially (7.21) is still valid, so  $\mathbf{res} \notin \mathbb{U}$ , and (7.24) and Lemma 2.5 finish the proof.  $\square$

With (7.19) we can also estimate the treatable vector length  $n$  in Algorithm 4.8 (**Transform**) depending on the number of extra extractions. The results are summarized in Table 7.1.

As before an algorithm with  $K$ -fold faithful rounding for input vectors of huge length may be developed as well.

TABLE 7.1  
*Maximum treatable dimension  $n$  depending on  $K$*

| $K$ | single precision | double precision    |
|-----|------------------|---------------------|
| 3   | 4094             | $9.5 \cdot 10^7$    |
| 4   | 65634            | $4.3 \cdot 10^{10}$ |
| 5   | $2.6 \cdot 10^5$ | $9.2 \cdot 10^{11}$ |
| 6   | $6.0 \cdot 10^5$ | $5.8 \cdot 10^{12}$ |
| 7   | $1.0 \cdot 10^6$ | $2.0 \cdot 10^{13}$ |

TABLE 9.1  
*Ratio of computing time for Dekker's algorithm 2.1 with branch to Knuth's algorithm 2.1*

| $n$       | CPU: Pentium IV 2.53GHz,<br>Compiler: Intel Fortran 8.1 | CPU: Pentium M 1GHz,<br>Compiler: Intel Fortran 8.1 |
|-----------|---|---|
| 100       | 1.662   | 6.571   |
| 1,000     | 2.247   | 8.429   |
| 10,000    | 2.532   | 8.700   |
| 100,000   | 2.538   | 2.313   |
| 1,000,000 | 2.462   | 1.644   |

**8. Dot products.** We briefly mention how to compute a faithful rounding of the value of a dot product of floating-point vectors, as well the rounded to nearest result. We already mentioned that for given  $a, b \in \mathbb{F}$ , Algorithm 1.1 (`TwoProduct`) together with Algorithm 1.2 (`Split`) compute  $x, y \in \mathbb{F}$  with

$$x + y = a \cdot b ,$$

provided no overflow and no underflow occurred. It might happen that an over- or underflow cancels, so that the exact result of a dot product is within the floating-point range. Although this occurs most likely only in constructed examples, it is easy to take care of.

For given vectors  $u, v \in \mathbb{F}^n$ , first split the index set  $\{1, \dots, n\}$  into large, small and medium, i.e. those indices where  $u_i \cdot v_i$  is near or in the overflow range, near or in the underflow range, and neither of those, respectively. This can of course be decided without causing an exception.

Then first handle the “large” indices with properly scaled input data, ignoring the rest. Using our estimations it can be decided whether the final result can possibly be influenced by the medium range. If this is not so, we are finished, possibly reporting overflow.

If not, treat the medium range the same way. Again it can be decided whether the “small” indices need to be treated. If this is the case, we proceed the same way with the scaled input data.

**9. Computational results.** In the following we give some computational results on different architectures and using different compilers. All programming and measurement was done by the second author.

In Section 2 we mentioned the severe slow down by branches due to lack of optimization. For example, Knuth's Algorithm 2.1 (`TwoSum`) may be replaced by Dekker's Algorithm 2.7 equipped with a branch  $|a| \geq |b|$ . We measured the elapsed time for both approaches and list the ratio of computing time in Table 9.1. For example, the ratio 1.662 means that Dekker's algorithm with 3 floating-point operations and one branch requires 66% more computing time than Knuth's algorithm with 6 floating operations. So Table 9.1 demonstrates how branches can ruin computing times. Fortunately we have ensured in all our error-free transformations of the sum of two floating-point numbers that Dekker's Algorithm 2.7 can be used *without* branch.

TABLE 9.2

*Ratio of computing times for extraction (CPU: Pentium IV 2.53GHz, Compiler: Intel Fortran 8.1)*

| $n$     | ExtractScalar | Dint | Dnint | Int=Dble | M/E   |
|---------|---------------|------|-------|----------|-------|
| 100     | 1.0           | 5.11 | 7.00  | 9.78     | 48.44 |
| 1,000   | 1.0           | 5.25 | 8.63  | 10.50    | 54.63 |
| 10,000  | 1.0           | 4.67 | 8.11  | 9.56     | 48.33 |
| 100,000 | 1.0           | 0.91 | 1.01  | 1.02     | 4.78  |

TABLE 9.3

*Ratio of computing times for extraction (CPU: Pentium M 1GHz, Compiler: Intel Fortran 8.1)*

| $n$     | ExtractScalar | Dint | Dnint | Int=Dble | M/E   |
|---------|---------------|------|-------|----------|-------|
| 100     | 1.0           | 4.45 | 6.00  | 17.94    | 31.97 |
| 1,000   | 1.0           | 3.79 | 5.03  | 14.84    | 26.26 |
| 10,000  | 1.0           | 3.22 | 4.24  | 12.30    | 21.65 |
| 100,000 | 1.0           | 1.39 | 1.75  | 4.82     | 8.43  |

All our algorithms are based on extractions, the split of a floating-point number with respect to  $\sigma$ , a power of 2 corresponding to a certain exponent. Since this operation is in the inner loop of all our algorithms, we payed special attention to this and designed Algorithm 3.1 (**ExtractScalar**) to be as fast as possible. This algorithm requires 3 floating-point operations and has no branch.

Another possibility to extract bits of a floating-point number  $p$  is proper scaling and rounding to integer. Some compilers offer two possibilities of such rounding, namely chopping and rounding to the nearest integer. In the Tables 9.2 and 9.3 the columns “Dint” and “Dnint” refer to those roundings, respectively. Another possibility of rounding a floating-point number is the assignment to a variable of type integer. One obstacle might be that an integer format with sufficient precision is not available. This approach is referred to by “Int=Dble”. Finally, the desired bits may be extracted by access to mantissa and exponent and some bit manipulation. This is referred to by “M/E”.

As can be seen from the Tables 9.2 to 9.3 our Algorithm 3.1 (**ExtractScalar**) is much faster than the other possibilities. Our algorithms directly benefit from this. There is a certain drop in the ratio for large dimension which is related to the cache size. We programmed all tests straightforwardly without blocking, so not too much attention must be paid to the last lines of the Tables 9.2 and 9.3 for huge dimension.

We designed our algorithms to take advantage of compiler optimization. However, sometimes we have to make sure that the latter is not overdone. The first line  $q = \text{fl}((\sigma + p) - \sigma)$  in Algorithm 3.1 (**ExtractScalar**), for example, may be erroneously optimized into  $q = p$ . This can, of course, be avoided by setting appropriate compiler options; however, this may slow down the whole computation. In this specific case the second author suggested a simple trick to overcome this by using  $q = \text{fl}(|\sigma + p| - \sigma)$  instead. This does not change the intended result since  $|p| \leq \sigma$  is assumed in the analysis (Lemma 3.2), it avoids unintended compiler optimization and it does not slow down the computation.

Next we tested our summation algorithms. For huge condition numbers larger than  $\text{eps}^{-1}$  we used our Algorithm 6.1 in [36] and residuals for the matrices from [42], where methods to generate arbitrarily ill-conditioned matrices exactly representable in floating-point are described. Dot products are transformed into sums by means of Algorithms 1.2 (**Split**) and 1.1 (**TwoProduct**).

We display results for dimension  $n = 1000$ . For other dimensions the results are similar, and again we observe a drop in the factors for huge dimensions due to limited cache size. In the following Tables 9.4 and 9.5 the column “DSum” refers to the ordinary recursive summation, the computing time of which is normed to 1.

TABLE 9.4

Ratio of computing times for summation ( $n = 1000$ , CPU: Pentium IV 2.53GHz, Compiler: Intel Visual Fortran 8.1)

| cond  | Dsum | AccSum      | Dint   | Dnint  | Int=Dble | Sum2         | Sum3         | XBLAS        |
|-------|------|-------------|--------|--------|----------|--------------|--------------|--------------|
| 1e8   | 1.0  | <b>7.00</b> | 21.80  | 28.20  | 36.20    | <b>18.20</b> | <b>32.60</b> | <b>25.80</b> |
| 1e16  | 1.0  | <b>7.02</b> | 21.83  | 28.20  | 36.20    | <b>18.20</b> | <b>32.60</b> | <b>25.80</b> |
| 1e32  | 1.0  | <b>9.60</b> | 37.00  | 44.80  | 62.60    | 18.20        | <b>32.60</b> | 25.80        |
| 1e64  | 1.0  | 14.17       | 58.17  | 77.67  | 101.33   | 18.20        | 32.60        | 25.80        |
| 1e128 | 1.0  | 29.00       | 124.20 | 161.00 | 215.40   | 18.20        | 32.60        | 25.80        |

TABLE 9.5

Ratio of computing times for summation ( $n = 1000$ , CPU: Pentium M 1GHz, Compiler: Intel Visual Fortran 8.1)

| cond  | Dsum | AccSum       | Dint   | Dnint  | Int=Dble | Sum2        | Sum3         | XBLAS        |
|-------|------|--------------|--------|--------|----------|-------------|--------------|--------------|
| 1e8   | 1.0  | <b>8.38</b>  | 23.75  | 29.81  | 75.56    | <b>8.59</b> | <b>15.41</b> | <b>14.35</b> |
| 1e16  | 1.0  | <b>8.40</b>  | 23.78  | 29.82  | 75.56    | <b>8.59</b> | <b>15.41</b> | <b>14.35</b> |
| 1e32  | 1.0  | <b>11.44</b> | 38.06  | 47.31  | 122.06   | 8.59        | <b>15.41</b> | 14.35        |
| 1e64  | 1.0  | 20.50        | 72.69  | 98.50  | 240.62   | 8.59        | 15.41        | 14.35        |
| 1e128 | 1.0  | 33.65        | 119.47 | 160.18 | 404.54   | 8.59        | 15.41        | 14.35        |

Column “AccSum” shows the computing time for our Algorithm 4.4 (AccSum), and the following columns “Dint”, “Dnint” and “Int=Dble” refer to AccSum with only the extraction changed as indicated.

Furthermore, columns “Sum2” and “Sum3” display the timing for our Algorithms 4.4 (Sum2) and Algorithm 4.8 (SumK) for  $K = 3$  out of [36]. These algorithms produce a result of quality as if computed in 2-fold or 3-fold precision, respectively. That means, for condition numbers up to  $\mathbf{eps}^{-1}$  or  $\mathbf{eps}^{-2}$ , respectively, we can expect a result accurate to the last bit. Finally, column “XBLAS” displays the computing time for the corresponding algorithm BLAS\_dsum\_x in XBLAS [31], the result of which is of the same quality as that of Sum2, namely as if computed in 2-fold precision.

For condition numbers up to  $\mathbf{eps}^{-1}$  we achieve a remarkable factor of about 7 or 8 compared to recursive summation. This is achieved by straightforward Fortran code without further optimization and without blocking.

For larger condition numbers gradually more extractions are needed, and correspondingly our Algorithm 4.4 (AccSum) requires more computing time. The timing is compared with ordinary (recursive) summation, the time of which is normed to 1.0, but the results of recursive summation are meaningless for condition numbers beyond  $\mathbf{eps}^{-1}$ . Note that the factors between the different options for extraction are similar to the ones obtained in Tables 9.2 and 9.3.

Similarly, the results of algorithm Sum2 and the results of XBLAS summation are meaningless for condition numbers beyond  $\mathbf{eps}^{-2}$ , and the results of algorithm Sum3 are meaningless for condition numbers beyond  $\mathbf{eps}^{-3}$ .

As we showed, algorithm AccSum always delivers a faithfully rounded result, i.e. accurate to the last bit. In those rows of Tables 9.4 and 9.5, where we can also expect algorithms Sum2, Sum3 and XBLAS to compute results accurate to the last bit, we displayed the corresponding computing times in **bold**. So the timings displayed in bold in one row refer to algorithms the results of which are of comparable accuracy.

As we see, our Algorithm 4.4 (AccSum) is always faster, sometimes much faster than the other approaches. The striking advantage of AccSum is that it guarantees a result *accurate* to the last bit rather than a result as if computed in some specified *precision*. Moreover, it only requires more computing time if necessary because of ill condition, and, in contrast to the other approaches, AccSum is perfectly suited for parallelization.

The computational results for dot products are very similar to the ones for summation and are therefore omitted.

Finally that we tried to find examples where the result of Algorithm 4.4 (**AccSum**) is not rounded to nearest. We treated dimensions from 10 to  $10^5$  and condition numbers as above. We especially used dimensions  $n = 2^M - 2$ . Fortunately, we have Algorithm 6.1 (**NearSum**) for reference. It is not so easy to find a long precision package delivering results always rounded to nearest; especially we observed problems with rounding tie to even. In several billion test cases we did not find one example with the result of **AccSum** not being the nearest floating-point number to the exact result.

**10. Appendix.** Following we give the proof of (5.10), that is

$$(10.1) \quad \tilde{R} := \tau_2 - \Delta \in \mathbb{F}$$

under the assumptions of Lemma 5.5 and for  $\sigma > \frac{1}{2}\mathbf{eps}^{-1}\mathbf{eta}$ . This implies  $R = \text{fl}(\tau_2 - \Delta) = \tau_2 - \Delta$ . Since  $\sigma$  is a power of 2, we have

$$(10.2) \quad \sigma \geq \mathbf{eps}^{-1}\mathbf{eta} .$$

We already know

$$(10.3) \quad \Delta = \mathbf{res} - \tau_1 \in \mathbb{F} .$$

We first prove some facts, namely

$$(10.4) \quad |\mathbf{res}| \geq \frac{3}{4}|\tau_1| ,$$

$$(10.5) \quad \frac{1}{2}\mathbf{eps}^{-1}\sigma \leq |\tau_1| \Rightarrow \mathbf{res} \in \{\text{pred}(\tau_1), \tau_1, \text{succ}(\tau_1)\} ,$$

$$(10.6) \quad \frac{1}{2}\mathbf{eps}^{-1}\sigma \leq |\tau_1| \text{ and } \mathbf{res} \neq \tau_1 \Rightarrow \text{sign}(\tau_2)\text{sign}(\Delta) \geq 0 .$$

Using (4.19), (4.11) and (4.10) it follows

$$|\mathbf{res}| \geq (1 - 2\mathbf{eps})|\tau_1| - n\mathbf{eps}\sigma \geq (1 - 2\mathbf{eps})|\tau_1| - 2^{-M-1}\text{ufp}(\tau_1) \geq \frac{3}{4}|\tau_1|$$

and prove (10.4). Furthermore, (5.8), (4.14), (4.11), (4.8) and the assumption  $\frac{1}{2}\mathbf{eps}^{-1}\sigma \leq |\tau_1|$  imply

$$(10.7) \quad \begin{aligned} |\tau'_2| &\leq (1 + \mathbf{eps})|\tau_2 + \tau_3| \leq |\tau_2| + \mathbf{eps}|\tau_2| + 2^M\mathbf{eps}\sigma \\ &\leq |\tau_2| + (\mathbf{eps}^2 + 2^{M+2}\mathbf{eps}^2)\text{ufp}(\tau_1) \\ &< |\tau_2| + \mathbf{eps} \cdot \text{ufp}(\tau_1) , \end{aligned}$$

so that  $\tau_1 = \text{fl}(\tau_1 + \tau_2)$ ,  $\mathbf{res} = \text{fl}(\tau_1 + \tau'_2)$  and (2.27) imply (10.5). If  $\mathbf{res} \neq \tau_1$ , then  $\text{sign}(\tau_2) = \text{sign}(\tau'_2) = \text{sign}(\mathbf{res} - \tau_1) = \text{sign}(\Delta)$  by (2.27) and the monotonicity of rounding, so (10.6).

**PROOF OF (5.10).** We can use the analysis of Algorithm 4.1 (**Transform**) and the notation in the proof of Lemma 4.3, especially (4.12). We distinguish several cases.

First, assume  $|\tau_1| < \sigma$ . Then (4.8) yields  $\tau_2 \in \mathbf{eps}\sigma\mathbb{Z}$  and  $|\tau_2| \leq \mathbf{eps} \cdot \text{ufp}(\tau_1) < \mathbf{eps}\sigma$ , so  $\tau_2 = 0$  and  $\tilde{R} = -\Delta \in \mathbb{F}$ .

Second, assume  $\sigma \leq |\tau_1| < \frac{3}{4}\mathbf{eps}^{-1}\sigma$ . Then  $\mathbf{res} = \text{fl}(\tau_1 + \tau'_2)$ , (2.17) and (2.11) yield  $\mathbf{res} \in \mathbf{eps} \cdot \text{ufp}(\tau_1) \subseteq \mathbf{eps}\sigma\mathbb{Z}$ , so  $\mathbf{res}, \tau_1, \tau_2 \in \mathbf{eps}\sigma\mathbb{Z}$  from (4.8) and  $\Delta, \tilde{R} \in \mathbf{eps}\sigma\mathbb{Z}$ . Moreover, (4.12), (2.20), (4.14) and (4.11) imply

$$(10.8) \quad \begin{aligned} |\tilde{R}| = |\tau_2 - \Delta| &= |\tau_2 - \tau'_2 + \delta_1| \leq |\tau_2 - \tau'_2| + \mathbf{eps}|\tau_1 + \tau'_2| \\ &\leq |\tau_3| + \mathbf{eps}|\tau_2 + \tau_3| + \mathbf{eps}|\tau_1| + \mathbf{eps}(1 + \mathbf{eps})|\tau_2 + \tau_3| \\ &\leq \mathbf{eps}|\tau_1| + \mathbf{eps}(2 + \mathbf{eps})|\tau_2| + (1 + \mathbf{eps})^2|\tau_3| \\ &\leq \mathbf{eps}(1 + 3\mathbf{eps})|\tau_1| + 2^M\mathbf{eps}\sigma \\ &< \sigma , \end{aligned}$$

and (2.14) and (10.2) prove  $\tilde{R} \in \mathbb{F}$ .

Third, assume  $\frac{3}{4}\mathbf{eps}^{-1}\sigma \leq |\tau_1| < \mathbf{eps}^{-1}\sigma$ . Then (10.5) implies  $\mathbf{res} \in \{\text{pred}(\tau_1), \tau_1, \text{succ}(\tau_1)\}$ . If  $\mathbf{res} = \tau_1$ , then  $\tilde{R} = \tau_2 \in \mathbb{F}$ . If  $\mathbf{res} \neq \tau_1$ , the assumption on  $\tau_1$  implies  $\text{ufp}(\tau_1) = \frac{1}{2}\mathbf{eps}^{-1}\sigma$  and

$$|\Delta| = |\mathbf{res} - \tau_1| = 2\mathbf{eps} \cdot \text{ufp}(\tau_1) = \sigma$$

by Lemma 2.2. But  $|\tau_2| \leq \mathbf{eps} \cdot \text{ufp}(\tau_1) = \frac{1}{2}\sigma$  by (4.8) and  $\text{sign}(\tau_2)\text{sign}(\Delta) \geq 0$  by (10.6) implies

$$\frac{1}{2}\sigma \leq |\tau_2 - \Delta| = |\tilde{R}| \leq \sigma \quad \text{and} \quad \tilde{R} \in \mathbf{eps}\sigma\mathbb{Z},$$

so (2.14) and (10.3) prove  $\tilde{R} \in \mathbb{F}$  also for that case.

Fourth, assume  $\mathbf{eps}^{-1}\sigma \leq |\tau_1| < 2\mathbf{eps}^{-1}\sigma$ . Then (10.4) yields

$$|\mathbf{res}| \geq \frac{3}{4}|\tau_1| > \frac{1}{2}\mathbf{eps}^{-1}\sigma,$$

so the distance of  $\mathbf{res}$  to its neighbors is at least  $2\mathbf{eps} \cdot \text{ufp}(\mathbf{res}) \geq \sigma$ . Using  $s + \varrho = \mathbf{res} + \delta$  as in (4.13) and (4.21) implies

$$|\varrho - \mathbf{res}| \leq |\delta| + |s| < \frac{1}{2}\mathbf{eps} \cdot \text{ufp}(\tau_1) + 2^{-M}\sigma = \left(\frac{1}{2} + 2^{-M}\right)\sigma < \sigma,$$

so  $\varrho, \mathbf{res} \in \mathbb{F}$  implies  $\varrho = \mathbf{res}$ . With (4.7) this yields

$$\begin{aligned} |\mathbf{res} - \tau_1| &\leq |\delta| + |\tau_2| + \left| \sum p'_i \right| < \left(\frac{1}{2}\mathbf{eps} + \mathbf{eps}\right)\text{ufp}(\tau_1) + 2^M\mathbf{eps}\sigma \\ &= \left(1 + \frac{1}{2} + 2^M\mathbf{eps}\right)\sigma < 2\sigma = 2\mathbf{eps} \cdot \text{ufp}(\tau_1). \end{aligned}$$

If  $\mathbf{res} = \tau_1$ , then  $\tilde{R} = \tau_2 \in \mathbb{F}$ . Otherwise, if  $\mathbf{res} \neq \tau_1$  and as by Lemma 2.2, the only possibility that the distance of the floating-point number  $\tau_1$  to another floating-point number  $\mathbf{res}$  is strictly less than  $2\mathbf{eps} \cdot \text{ufp}(\tau_1)$  is  $|\tau_1| = \text{ufp}(\tau_1) = \mathbf{eps}^{-1}\sigma$  and

$$|\Delta| = |\mathbf{res} - \tau_1| = \mathbf{eps} \cdot \text{ufp}(\tau_1) = \sigma.$$

Moreover, (10.5) implies  $\text{sign}(\tau_2)\text{sign}(\Delta) \geq 0$ , so that with  $\tilde{R} = \tau_2 - \Delta$  with  $|\tau_2| \leq \mathbf{eps} \cdot \text{ufp}(\tau_1) = \sigma$  we have

$$0 \leq |\tau_2 - \Delta| = |\tilde{R}| \leq \sigma \quad \text{and} \quad \tilde{R} \in \mathbf{eps}\sigma\mathbb{Z},$$

and again (2.14) and (10.3) prove  $\tilde{R} \in \mathbb{F}$  also for that case.

Fifth and last, assume  $|\tau_1| \geq 2\mathbf{eps}^{-1}\sigma$ . By (4.3) and (4.2) we have for the final value  $m$  in Algorithm 4.1 (**Transform**)

$$|t^{(m-1)}| = |\tau_1 + \tau_2 - \tau^{(m)}| \geq (1 - \mathbf{eps})|\tau_1| - \sigma > \mathbf{eps}^{-1}\sigma.$$

Suppose  $m \geq 2$ , then the “until”-condition in **Transform** implies

$$|t^{(m-1)}| < 2^{2M+1}\mathbf{eps}\sigma_{m-2} = 2^{M+1}\sigma_{m-1} = 2^{M+1}\sigma < \mathbf{eps}^{-1}\sigma,$$

a contradiction. Thus  $m = 1$ ,  $t^{(m-1)} = t^{(0)} = \varrho$ ,  $\tau_1 + \tau_2 = \varrho + \tau^{(m)}$  and  $\mathbf{res} = \text{fl}(\varrho + \tau^{(m)})$ . But (10.4) yields

$$|\mathbf{res}| \geq \frac{3}{4}|\tau_1| > \mathbf{eps}^{-1}\sigma,$$

and the distance of  $\mathbf{res}$  to its neighbors is at least  $2\mathbf{eps} \cdot \text{ufp}(\mathbf{res}) \geq 2\sigma$ . Hence  $|\tau^{(m)}| < \sigma_{m-1} = \sigma$  implies

$$\varrho = \text{fl}(\varrho + \tau^{(m)}) = \text{fl}(\tau_1 + \tau_2) = \tau_1,$$

it follows  $\Delta = 0$  and  $\tilde{R} = \tau_2 \in \mathbb{F}$  and also this case is finished.

All together it shows  $R = \text{fl}(\tau_2 - \Delta) = \text{fl}(\tilde{R}) = \tilde{R} = \tau_2 - \Delta$ . This closes the proof of (5.10).  $\square$

Following is executable Matlab code for Algorithm 4.4 including acceleration for zero sums and elimination of zero summands for that case. Moreover, the algorithms `ExtractVector` and `FastTwoSum` are expanded. Note that the Matlab function `nextpow2(f)` returns the smallest  $k$  such that  $2^k \geq |f|$ , while Algorithm 3.5 (`NextPowerTwo`) returns  $2^k$ . Accordingly, the variable `Ms` refers to  $2^M$  in Algorithm 4.4 (`AccSum`).

ALGORITHM 10.1. *Executable Matlab code for Algorithm 4.4 (AccSum) including check for zero sum.*

```
function res = AccSum(p)
% For given vector p, result res is the exact sum of p_i faithfully rounded.
% Acceleration for zero sums is included.
%
n = length(p);                % initialization
mu = max(abs(p));             % abs(p_i) <= mu
if ( n==0 ) | ( mu==0 )      % no or only zero summands
    res = 0;
    return
end
Ms = NextPowerTwo(n+2);       % n+2 <= 2^M
sigma = Ms*NextPowerTwo(mu);  % first extraction unit
phi = 2^(-53)*Ms;            % factor to decrease sigma
factor = 2^(-52)*Ms*Ms;      % factor for sigma check
%
t = 0;
while 1
    q = ( sigma + p ) - sigma; % [tau,p] = ExtractVector(sigma,p);
    tau = sum(q);              % sum of leading terms
    p = p - q;                 % remaining terms
    tau1 = t + tau;            % new approximation
    if ( abs(tau1)>=factor*sigma ) | ( sigma<=realmin )
        tau2 = tau - ( tau1 - t ); % [tau1,tau2] = FastTwoSum(t,tau)
        res = tau1 + ( tau2 + sum(p) ); % faithfully rounded final result
        return
    end
    t = tau1;                  % sum t+tau exact
    if t==0                    % accelerate case sum(p)=0
        res = AccSum(p(p~=0)); % recursive call, zeros eliminated
        return
    end
    sigma = phi*sigma;         % new extraction unit
end
```

**11. Summary.** We presented algorithms for summation with provably faithfully rounded result and provably rounded to nearest result in working precision. Furthermore, an algorithm was presented for summation with faithfully rounded result in  $K$ -fold accuracy. The paper contains as well the ingredients to compute a rounded to nearest result in  $K$ -fold accuracy. All our algorithms use only ordinary floating-point addition, subtraction and multiplication, no branches in the inner loops and no special operations. Similar

algorithms for dot products are sketched.

We showed that our summation Algorithm 4.4 (AccSum) is always faster, sometimes much faster than other approaches. The striking advantage of AccSum is that it guarantees a result *accurate* to the last bit rather than a result as if computed in some specified *precision*. Moreover, in contrast to other approaches, AccSum is perfectly suited for parallelization.

The algorithms are based on so-called error-free transformations. We hope to see these computationally and mathematically highly interesting operations in future computer architectures and floating-point standards.

**Acknowledgement.** The first author wishes to express his thankfulness that the paper could be written during stays at Waseda University supported by the Grant-in-Aid for Specially Promoted Research from the Mext, Japan. The first author also wishes to thank his students of the winter term 2004/05 and the summer term 2005 for their patience and constructive comments. The second author would like to express his sincere thanks to Dr. Y. Ushiro for his stimulating discussions.

#### REFERENCES

- [1] *MPFR: A C library for multiple-precision floating-point computations with exact rounding*. available at <http://www.mpfr.org>.
- [2] F. AVNAIM, J. BOISSONNAT, O. DEVILLERS, F. PREPARATA, AND M. YVINEC, *Evaluating signs of determinants using single-precision arithmetic*, *Algorithmica*, 17 (1997), pp. 111–132.
- [3] D. BAILEY, *A Fortran-90 based multiprecision system*, *ACM Trans. Math. Software*, 21 (1995), pp. 379–387.
- [4] G. BOHLENDER, *Floating-point computation of functions with maximum accuracy*, *IEEE Trans. Comput.*, C-26 (1977), pp. 621–632.
- [5] R. BRENT, *A Fortran Multiple-Precision Arithmetic Package*, tech. report, Dept. of Computer Science, Australian National University, Canberra, 1975.
- [6] H. BRÖNNIMANN, C. BURNIKEL, AND S. PION, *Interval arithmetic yields efficient dynamic filters for computational geometry*, *Discrete Applied Mathematics*, 109 (2001), pp. 25–47.
- [7] H. BRÖNNIMANN AND M. YVINEC, *Efficient exact evaluation of signs of determinants*, *Algorithmica*, 27 (2000), pp. 21–56.
- [8] K. CLARKSON, *Safe and effective determinant evaluation*, in *Proceedings of the 33rd Annual Symposium on Foundations of Computer Science (Pittsburgh, PA, IEEE Computer Society Press, 1992*, pp. 387–395.
- [9] M. DAUMAS AND D. MATULA, *Validated roundings of dot products by sticky accumulation*, *IEEE Trans. Computers*, 46 (1997), pp. 623–629.
- [10] T. J. DEKKER, *A floating-point technique for extending the available precision*, *Numer. Math.*, 18 (1971), pp. 224–242.
- [11] J. DEMMEL AND Y. HIDA, *Accurate and Efficient Floating Point Summation*, *SIAM J. Sci. Comput.*, 25 (2003), pp. 1214–1248.
- [12] J. DEMMEL AND Y. HIDA, *Fast and accurate floating point summation with application to computational geometry*, *Numerical Algorithms*, 37 (2004), pp. 101–112.
- [13] *GNU multiple precision arithmetic library (GMP), version 4.1.2*, 2003. Code and documentation available at <http://swox.com/gmp>.
- [14] S. GRAILLAT, P. LANGLOIS, AND N. LOUVET, *Compensated Horner Scheme*, tech. report, Laboratoire LP2A, University of Perpignan, 2005.
- [15] J. HAUSER, *Handling floating-point exceptions in numeric programs*, *ACM Trans. Program. Lang. Syst.*, 18 (1996), pp. 139–174.
- [16] C. HECKER, *Let's get to the (floating) point*, *Game Developer*, 2 (1996), pp. 19–24.
- [17] N. HIGHAM, *The accuracy of floating point summation*, *SIAM J. Sci. Comput.*, 14 (1993), pp. 783–799.
- [18] ———, *Accuracy and Stability of Numerical Algorithms*, SIAM Publications, Philadelphia, 2nd ed., 2002.
- [19] C. HOFFMANN, *Robustness in geometric computations*, *JCISE*, 1 (2001), pp. 143–156.
- [20] ANSI/IEEE, *IEEE Standard for Binary Floating Point Arithmetic*, IEEE, New York, Std 754–1985 ed., 1985.
- [21] M. JANKOWSKI, A. SMOKTUNOWICZ, AND WOŹNIAKOWSKI, *A note on floating-point summation of very many terms*, *J. Information Processing and Cybernetics-EIK*, 19 (1983), pp. 435–440.
- [22] M. JANKOWSKI AND WOŹNIAKOWSKI, *The accurate solution of certain continuous problems using only single precision arithmetic*, *BIT*, 25 (1985), pp. 635–651.
- [23] W. KAHAN, *A survey of error analysis*, in *Proceedings of the IFIP Congress, Information Processing 71, North-Holland, Amsterdam, The Netherlands, 1972*, pp. 1214–1239.
- [24] ———, *Implementation of Algorithms (Lecture Notes by W.S. Haugeland and D. Hough)*, Tech. Report 20, Department of Computer Science, University of California, Berkeley, CA, 1973.

- [25] A. KIELBASZIŃSKI, *Summation algorithm with corrections and some of its applications*, Math. Stos, 1 (1973), pp. 22–41.
- [26] D. KNUTH, *The Art of Computer Programming: Seminumerical Algorithms*, vol. 2, Addison Wesley, Reading, Massachusetts, 1969.
- [27] S. KRISHNAN, M. FOSKEY, T. CULVER, J. KEYSER, AND D. MANOCHA, *PRECISE: efficient multiprecision evaluation of algebraic roots and predicates for reliable geometric computation*, in SCG '01: Proceedings of the 17th annual symposium on Computational geometry, New York, NY, USA, 2001, ACM Press, pp. 274–283.
- [28] U. KULISCH AND W. L. MIRANKER, *The arithmetic of the digital computer: A new approach*, SIAM Review, 28 (1986), pp. 1–40.
- [29] P. LANGLOIS AND N. LOUVET, *Solving Trinagular Systems More Accurately and Efficiently*, tech. report, Laboratoire LP2A, University of Perpignan, 2005.
- [30] H. LEUPRECHT AND W. OBERAIGNER, *Parallel algorithms for the rounding exact summation of floating point numbers*, Computing, 28 (1982), pp. 89–104.
- [31] X. LI, J. DEMMEL, D. BAILEY, G. HENRY, Y. HIDA, J. ISKANDAR, W. KAHAN, S. KANG, A. KAPUR, M. MARTIN, B. THOMPSON, T. TUNG, AND D. YOO, *Design, implementation and testing of extended and mixed precision BLAS*, ACM Trans. Math. Softw., 28 (2002), pp. 152–205.
- [32] S. LINNAINMAA, *Software for doubled-precision floating point computations*, ACM Trans. Math. Soft., 7 (1981), pp. 272–283.
- [33] M. MALCOLM, *On accurate floating-point summation*, Comm. ACM, 14 (1971), pp. 731–736.
- [34] O. MØLLER, *Quasi double precision in floating-point arithmetic*, BIT, 5 (1965), pp. 37–50.
- [35] A. NEUMAIER, *Rundungsfehleranalyse einiger Verfahren zur Summation endlicher Summen*, ZAMM, 54 (1974), pp. 39–51.
- [36] T. OGITA, S. RUMP, AND S. OISHI, *Accurate Sum and Dot Product*, SIAM J. Sci. Comput., 26 (2005), pp. 1955–1988.
- [37] M. PICHAT, *Correction d'une somme en arithmétique à virgule flottante*, Numer. Math., 19 (1972), pp. 400–406.
- [38] D. PRIEST, *Algorithms for arbitrary precision floating point arithmetic*, in Proceedings of the 10th Symposium on Computer Arithmetic, P. Kornerup and D. Matula, eds., Grenoble, France, 1991, IEEE Computer Society Press, pp. 132–145.
- [39] ———, *On properties of floating point arithmetics: Numerical stability and the cost of accurate computations*, PhD thesis, Mathematics Department, University of California at Berkeley, CA, 1992.
- [40] D. ROSS, *Reducing truncation errors using cascading accumulators*, Communication of ACM, 8 (1965), pp. 32–33.
- [41] S. RUMP, *Kleine Fehlerschranken bei Matrixproblemen*, PhD thesis, Universität Karlsruhe, 1980.
- [42] ———, *A class of arbitrarily ill-conditioned floating-point matrices*, SIAM J. Mat. Anal. Appl., 12 (1991), pp. 645–653.
- [43] S. SCHIRRA, *Precision and robustness in geometric computations*, in Algorithmic Foundations of Geographic Information Systems, M. van Kreveld, J. Nievergelt, T. Roos, and P. Widmayer, eds., Lecture Notes in Computer Science 1340, Springer Verlag, Berlin, 1997, pp. 255–287.
- [44] J. SHEWCHUK, *Adaptive precision floating-point arithmetic and fast robust geometric predicates*, Discrete Comput. Geom., 18 (1997), pp. 305–363.
- [45] G. ZIELKE AND V. DRYGALLA, *Genauere Lösung linearer Gleichungssysteme*, GAMM Mitt. Ges. Angew. Math. Mech., 26 (2003), pp. 7–108.