# ACCURATE FLOATING-POINT SUMMATION PART II: SIGN, $K$-FOLD FAITHFUL AND ROUNDING TO NEAREST *

SIEGFRIED M. RUMP †, TAKESHI OGITA ‡, AND SHIN'ICHI OISHI §

**Abstract.** In this Part II of this paper we first refine the analysis of error-free vector transformations presented in Part I. Based on that we present an algorithm for calculating the rounded-to-nearest result of $s := \sum p_i$ for a given vector of floating-point numbers $p_i$, as well as algorithms for directed rounding. A special algorithm for computing the sign of $s$ is given, also working for huge dimensions. Assume a floating-point working precision with relative rounding error unit eps. We define and investigate a $K$-fold faithful rounding of a real number $r$. Basically the result is stored in a vector $\mathtt{Res}_\nu$ of $K$ non-overlapping floating-point numbers such that $\sum \mathtt{Res}_\nu$ approximates $r$ with relative accuracy $\mathtt{eps}^K$, and replacing $\mathtt{Res}_K$ by its floating-point neighbors in $\sum \mathtt{Res}_\nu$ forms a lower and upper bound for $r$. For a given vector of floating-point numbers with exact sum $s$, we present an algorithm for calculating a $K$-fold faithful rounding of $s$ using solely the working precision. Furthermore, an algorithm for calculating a faithfully rounded result of the sum of a vector of huge dimension is presented. Our algorithms are fast in terms of measured computing time because they allow good instruction-level parallelism, they neither require special operations such as access to mantissa or exponent, they contain no branch in the inner loop, nor do they require some extra precision: The only operations used are standard floating-point addition, subtraction and multiplication in one working precision, for example double precision. Certain constants used in the algorithms are proved to be optimal.

**1. Introduction, notation and basic facts.** We will present fast algorithms to compute approximations of high quality of the sum of a vector $p_i$ of floating-point numbers. Since sums of floating-point numbers are ubiquitous in scientific computations, there is a vast amount of literature to that; excellent surveys can be found in [9, 14].

In Part I [22] of this paper we gave a fast algorithm to compute a faithfully rounded result of $\sum p_i$. Our methods are based on *error-free transformations*. For example, Knuth [11] gave an algorithm (cf. Part I, Algorithm 2.1) transforming the sum $a + b$ of two floating-point numbers into a sum $x + y$, where $x$ is the usual floating-point approximation of the sum and $y$ comprises of the exact error. Surprisingly, $x$ and $y$ can be calculated using only 6 ordinary floating-point operations. Such error-free transformations receive interest in many areas [1, 8, 13, 14, 17, 18, 19, 20, 24, 25, 26, 27].

More background, an overview of existing methods and more details are given in Part I [22] of this paper.

This Part II of our paper extends the results of Part I [22] in various ways, and it is organized as follows. For the often delicate estimations we developed a framework for the analysis in Part I. Those techniques and main results are summarized in Section 2. Next we redefine the error-free vector transformation in Algorithm

†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).

‡Department of Mathematics, Tokyo Woman's Christian University, 2–6–1 Zempukuji, Suginami-ku, Tokyo 167-8585, Japan, and Visiting Associate Professor at Waseda University, Faculty of Science and Engineering, 3–4–1 Okubo, Shinjuku-ku, Tokyo 169–8555, Japan (ogita@lab.twcu.ac.jp).

§Department of Computer Science, Faculty of Science and Engineering, Waseda University, 3–4–1 Okubo, Shinjuku-ku, Tokyo 169–8555, Japan (oishi@waseda.jp).

3.3 (`Transform`) and refine the analysis of Part I. This is the key to all algorithms in Part II. We especially allow for huge vector lengths up to about $\texttt{eps}^{-1}$, and give as a first example in Section 4 an algorithm to compute the sign of $\sum p_i$. We show that the constant in the stopping criterion is best possible.

In Part I we developed an algorithm to compute a faithful rounding `res` of the sum $s$ of a vector of floating-point numbers. This means that there is no floating-point number between `res` and $s$. Especially, if $s$ is itself a floating-point number or is in the underflow range, then $\texttt{res} = s$. In the following Section 5 of this paper we define and investigate $K$-fold faithful rounding. Suppose a floating-point working precision with relative rounding error unit `eps` to be given. Then $K$-fold faithful rounding of $s := \sum p_i$ means basically that a vector $\texttt{Res}_\nu$ of $K$ floating-point numbers is computed such that $\sum \texttt{Res}_\nu$ is an approximation of relative accuracy $\texttt{eps}^K$, and $\texttt{Res}_K$ is a faithful rounding of $s - \sum_{\nu=1}^{K-1} \texttt{Res}_\nu$. This implies that replacing $\texttt{Res}_K$ by its two floating-point neighbors in $\sum \texttt{Res}_\nu$ produces a lower and upper bound for $s$. After developing the theoretical background in Section 5, we present a fast algorithm for computing a $K$-fold faithful rounding $\texttt{Res}_\nu$ of $\sum p_i$ in Section 6. Moreover we show that the sequence $\texttt{Res}_\nu$ is non-overlapping.

In the following Section 7 we develop a rounding-to-nearest algorithm for $s = \sum p_i$. This is the ultimate accuracy of an approximation of $s$ by a single floating-point number; however, it necessarily comes with a burden: To compute a faithful rounding of $s$, it suffices to know $s$ up to some error margin, whereas the rounded-to-nearest result may ultimately require to know $s$ exactly, namely if $s$ is the midpoint of two adjacent floating-point numbers.

Our Algorithm 4.5 (`AccSum`) presented in Part I for computing a faithfully rounded result of $\sum p_i$ has the charming property that the computing time is proportional to the logarithm of the condition number: The more difficult the problem is, the more computing time is needed. This is also true for our rounding-to-nearest algorithm, however, the "difficulty" depends on the nearness of the exact result $\sum p_i$ to the midpoint of two adjacent floating-point numbers. In Section 7 we also give algorithms for computing $\sum p_i$ with directed rounding.

For our Algorithm 4.5 (`AccSum`) presented in Part I the vector length was limited to about $\sqrt{\texttt{eps}^{-1}}$. In Section 8 we extend the range of applicability to vector lengths near $\texttt{eps}^{-1}$. We conclude the paper with computational results on a Pentium 4, Itanium 2 and Athlon 64 processor. For all algorithms presented in Part I and II of this paper and in [18] we put a Matlab reference code on `http://www.ti3.tu-harburg.de/rump` .

As in [18] and [22], all theorems, error analysis and proofs are due to the first author of the present paper.

**2. Notation and basic facts.** We use the notation and a number of results of Part I [22] of this paper. For convenience, some of the main results are summarized in the following, for more details, cf. [22].

The set of floating-point numbers is denoted by $\mathbb{F}$, and $\mathbb{U}$ denotes the set of subnormal 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 larger floating-point number, is denoted by `eps`, and the underflow unit by `eta`, that is the smallest positive (subnormal) floating-point number. For IEEE 754 double precision we have $\texttt{eps} = 2^{-53}$ and $\texttt{eta} = 2^{-1074}$. Then $\frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$ is the smallest positive normalized floating-point number, and for $f \in \mathbb{F}$ we have

$$(2.1) \qquad\qquad f \in \mathbb{U} \Leftrightarrow 0 \leq |f| \leq \frac{1}{2}\texttt{eps}^{-1}\texttt{eta} .$$

Note that for $f \in \mathbb{U}$, $f \pm \texttt{eta}$ are the floating-point neighbors of $f$. 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. An expression like $\text{fl}\big(\sum p_i\big)$ implies inherently that summation may be performed in any order. We assume floating-point operations in rounding to nearest corresponding to the IEEE 754 arithmetic standard [10].

$$\sigma := \mathrm{ufp}(f) \qquad\qquad 2\,\mathrm{eps}\,\sigma$$

FIG. 2.1. *The unit in the first place and unit in the last place of a normalized floating-point number*

In Part I we introduced the "unit in the first place" (ufp) or leading bit of a real number by

$$(2.2) \qquad\qquad 0 \neq r \in \mathbb{R} \quad \Rightarrow \quad \mathrm{ufp}(r) := 2^{\lfloor \log_2 |r| \rfloor} \ ,$$

where we set $\mathrm{ufp}(0) := 0$. This gives a convenient way to characterize the bits of a normalized floating-point number $f$: They range between the leading bit $\mathrm{ufp}(f)$ and the unit in the last place $2\mathrm{eps} \cdot \mathrm{ufp}(f)$. The situation is depicted in Figure 2.1.

In our analysis we will frequently view a floating-number as a scaled integer. For $\sigma = 2^k, k \in \mathbb{Z}$, we use the set $\mathrm{eps}\sigma\mathbb{Z}$, which can be interpreted as a set of fixed point numbers with smallest positive number $\mathrm{eps}\sigma$. Of course, $\mathbb{F} \subseteq \mathrm{eta}\mathbb{Z}$. These two concepts, the unit in the first place $\mathrm{ufp}(\cdot)$ together with $f \in \mathbb{F} \Rightarrow f \in 2\mathrm{eps} \cdot \mathrm{ufp}(f)\mathbb{Z}$ proved to be very useful in the often delicate analysis of our algorithms. Note that (2.2) is independent of some floating-point format and it applies to real numbers as well: $\mathrm{ufp}(r)$ is the value of the first nonzero bit in the binary representation of $r$. It follows

$$(2.3) \qquad\qquad 0 \neq r \in \mathbb{R} \quad \Rightarrow \quad \mathrm{ufp}(r) \leq |r| < 2\mathrm{ufp}(r)$$

$$(2.4) \qquad \sigma' = 2^m,\ m \in \mathbb{Z} \ \text{and} \ \sigma' \geq \sigma \quad \Rightarrow \quad \mathrm{eps}\sigma'\mathbb{Z} \subseteq \mathrm{eps}\sigma\mathbb{Z}$$

$$(2.5) \qquad r \in \mathrm{eps}\sigma\mathbb{Z},\ |r| \leq \sigma \ \text{and} \ \mathrm{eps}\sigma \geq \mathrm{eta} \quad \Rightarrow \quad r \in \mathbb{F}$$

$$(2.6) \qquad a, b \in \mathbb{F} \cap \mathrm{eps}\sigma\mathbb{Z} \ \text{and} \ \delta := \mathrm{fl}(a+b) - (a+b) \quad \Rightarrow \quad \mathrm{fl}(a+b),\ a+b,\ \delta \in \mathrm{eps}\sigma\mathbb{Z}$$

$$(2.7) \qquad n\mathrm{eps} \leq 1,\ a_i \in \mathbb{F} \ \text{and} \ |a_i| \leq \sigma \quad \Rightarrow \quad |\mathrm{fl}(\textstyle\sum_{i=1}^n a_i)| \leq n\sigma \ \text{and} \\ |\mathrm{fl}(\textstyle\sum_{i=1}^n a_i) - (\textstyle\sum_{i=1}^n a_i)| \leq \frac{n(n-1)}{2}\mathrm{eps}\sigma$$

$$(2.8) \qquad a, b \in \mathbb{F},\ a \neq 0 \quad \Rightarrow \quad \mathrm{fl}(a+b) \in \mathrm{eps} \cdot \mathrm{ufp}(a)\mathbb{Z} \ ,$$

see (2.9) through (2.20) in Part I of this paper. The fundamental error bound for floating-point addition is

$$(2.9) \qquad f = \mathrm{fl}(a+b) \quad \Rightarrow \quad f = a + b + \delta \ \text{with} \ |\delta| \leq \mathrm{eps} \cdot \mathrm{ufp}(a+b) \leq \mathrm{eps} \cdot \mathrm{ufp}(f) \leq \mathrm{eps}|f| \ ,$$

cf. (2.19) in part I. Note that this improves the standard error bound $\mathrm{fl}(a+b) = (a+b)(1+\varepsilon)$ for $a, b \in \mathbb{F}$ and $|\varepsilon| \leq \mathrm{eps}$ by up to a factor 2. In many estimations we desperately need this factor. Also note that (2.9) is also true in the underflow range, in fact addition (and subtraction) is exact if $\mathrm{fl}(a \pm b) \in \mathbb{U}$.

Next we note a useful sufficient criterion [cf. (2.21) in Part I] to decide that no error occurred in the sum of two floating-point numbers. For $a, b \in \mathbb{F}$ and $\sigma = 2^k,\ k \in \mathbb{Z}$,

$$(2.10) \qquad \begin{aligned} a, b \in \mathrm{eps}\sigma\mathbb{Z} \ \text{and} \ \mathrm{fl}(|a+b|) < \sigma \quad &\Rightarrow \quad \mathrm{fl}(a+b) = a+b \qquad \text{and} \\ a, b \in \mathrm{eps}\sigma\mathbb{Z} \ \text{and} \ |a+b| \leq \sigma \quad &\Rightarrow \quad \mathrm{fl}(a+b) = a+b \ . \end{aligned}$$

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

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

Using the ufp concept, the predecessor and successor of a floating-point number can be characterized as follows. Note that $0 \neq |f| = \mathrm{ufp}(f)$ is equivalent to $f$ being a power of 2.

LEMMA 2.1. *Let a floating-point number $0 \neq f \in \mathbb{F}$ be given. Then*

$$f \notin \mathbb{U} \ \text{and} \ |f| \neq \mathrm{ufp}(f) \quad \Rightarrow \quad \mathrm{pred}(f) = f - 2\mathrm{eps} \cdot \mathrm{ufp}(f) \ \text{and} \ f + 2\mathrm{eps} \cdot \mathrm{ufp}(f) = \mathrm{succ}(f) \ ,$$

$$f \notin \mathbb{U} \quad \text{and} \quad f = \text{ufp}(f) \quad \Rightarrow \quad \text{pred}(f) = (1 - \text{eps})f \quad \text{and} \quad (1 + 2\text{eps})f = \text{succ}(f) \ ,$$

$$f \notin \mathbb{U} \quad \text{and} \quad f = -\text{ufp}(f) \quad \Rightarrow \quad \text{pred}(f) = (1 + 2\text{eps})f \quad \text{and} \quad (1 - \text{eps})f = \text{succ}(f) \ ,$$

$$f \in \mathbb{U} \quad \Rightarrow \quad \text{pred}(f) = f - \texttt{eta} \quad \text{and} \quad f + \texttt{eta} = \text{succ}(f) \ .$$

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

$$(2.11) \qquad\qquad f - 2\text{eps} \cdot \text{ufp}(f) \leq \text{pred}(f) < \text{succ}(f) \leq f + 2\text{eps} \cdot \text{ufp}(f) \ .$$

A main concept of both parts of this paper is *faithful rounding* [5, 20, 4]. A floating-point number $f$ is called a faithful rounding of a real number $r$ if there is no other floating-point number between $f$ and $r$. It follows that $f = r$ in case $r \in \mathbb{F}$.

DEFINITION 2.2. *A floating-point number $f \in \mathbb{F}$ is called a* faithful rounding *of a real number $r \in \mathbb{R}$ if*

$$(2.12) \qquad\qquad\qquad\qquad \text{pred}(f) < r < \text{succ}(f) \ .$$

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

The following sufficient criterion was given in Lemma 2.4 in Part I. Note that for the computation of a faithful rounding of a real number $r$ we only need to know $r$ up to a certain error margin. The rounded-to-nearest $\text{fl}(r)$ computed by Algorithm 7.4, however, ultimately requires to know $r$ precisely, namely if $r$ is the midpoint between two adjacent floating-point numbers.

LEMMA 2.3. *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 $|\delta| < \frac{1}{2}\texttt{eta}$. Then $\tilde{r} \in \square(r + \delta)$, that means $\tilde{r}$ is a faithful rounding of $r + \delta$.*

A main principle of both parts of this paper is the error-free transformation of the sum of floating-point numbers. For the sum of two floating-point numbers Knuth [11] gave an algorithm transforming $a + b = x + y$ for general $a, b \in \mathbb{F}$ with $x = \text{fl}(a + b)$. It requires 6 floating-point operations. The following faster version by Dekker [5] applies if $a, b$ are somehow sorted.

ALGORITHM 2.4. *Compensated summation of two floating-point numbers.*

$$\text{function } [x, y] = \texttt{FastTwoSum}(a, b)$$
$$x = \text{fl}(a + b)$$
$$q = \text{fl}(x - a)$$
$$y = \text{fl}(b - q)$$

In Part I, Lemma 2.6 we analyzed the algorithm as follows.

LEMMA 2.5. *Let $a, b$ be floating-point numbers with $a \in 2\text{eps} \cdot \text{ufp}(b)\mathbb{Z}$. Let $x, y$ be the results produced by Algorithm 2.4 (`FastTwoSum`) applied to $a, b$. Then*

$$(2.13) \qquad\qquad x + y = a + b \ , \quad x = \text{fl}(a + b) \quad \text{and} \quad |y| \leq \text{eps} \cdot \text{ufp}(a + b) \leq \text{eps} \cdot \text{ufp}(x) \ .$$

*Furthermore,*

$$(2.14) \qquad\qquad\qquad q = \text{fl}(x - a) = x - a \quad \text{and} \quad y = \text{fl}(b - q) = b - q \ ,$$

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

**3. Error-free transformations.** The main principle of the error-free vector transformation in Part I is the extraction of a vector into a sum of higher order parts and a vector of lower order parts. The splitting is chosen so that the higher order parts add without error.

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

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

The following was proved in Part I, Theorem 3.4.

THEOREM 3.2. *Let $\tau$ and $p'$ be the results of Algorithm 3.1 (*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 $k \in \mathbb{Z}$, $n < 2^M$ for some $M \in \mathbb{N}$ and $|p_i| \leq 2^{-M}\sigma$ for all $i$. Then*

$$(3.1) \qquad \sum_{i=1}^{n} p_i = \tau + \sum_{i=1}^{n} p'_i , \quad \max|p'_i| \leq \texttt{eps}\sigma , \quad |\tau| \leq (1 - 2^{-M})\sigma < \sigma \quad and \quad \tau \in \texttt{eps}\sigma\mathbb{Z} .$$

Based on that we derived in Part I a transformation of a vector $p_i$ into an approximation of its sum and a remainder part, namely $\sum p_i = \tau_1 + \tau_2 + \sum p'_i$, an error-free transformation. Now we refine this algorithm and its analysis by introducing an offset, by a parameterized stopping criterion and by allowing for huge vector lengths.

ALGORITHM 3.3. *Transformation of a vector $p^{(0)}$ plus offset $\varrho$ (depending on parameter $\Phi$).*

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

Note that this is a preliminary version of Transform with all intermediate results uniquely identified using indices to ease readability and analysis. The original Algorithm 4.1 (Transform) in Part I was rewritten into the final version Algorithm 4.4 by omitting indices, overwriting results and an additional check for zero. Completely similarly Algorithm 3.3 is to be rewritten in an actual implementation.

LEMMA 3.4. *Let $\tau_1, \tau_2, p^{(m)}, \sigma$ be the results of Algorithm 3.3 (*Transform*) applied to a nonzero vector of floating-point numbers $p_i^{(0)}, 1 \leq i \leq n$, and let $\varrho \in \mathbb{F}$. Define $M := \lceil \log_2(n + 2) \rceil$ and assume $2^M \texttt{eps} < 1$. Furthermore, assume that $\varrho \in \texttt{eps}\sigma_0\mathbb{Z}$ is satisfied for $\mu := \max_i |p_i^{(0)}|$ and $\sigma_0 = 2^{M + \lceil \log_2 \mu \rceil}$. Assume the parameter $\Phi$ in the "until"-condition to be a power of 2 satisfying $\texttt{eps} \leq \Phi \leq 1$. Denote $s := \sum_{i=1}^{n} p_i^{(0)}$.*

*Then Algorithm 3.3 will stop, and*

$$(3.2) \qquad s + \varrho = t^{(m-1)} + \tau^{(m)} + \sum_{i=1}^{n} p_i^{(m)}$$

(3.3)    $\max |p_i^{(m)}| \leq \text{eps}\,\sigma_{m-1}$ ,    $|\tau^{(m)}| \leq (1 - 2^{-M})\sigma_{m-1} < \sigma_{m-1}$    and    $t^{(m-1)}, \tau^{(m)} \in \text{eps}\,\sigma_{m-1}\mathbb{Z}$

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

(3.4)    $\tau_1 + \tau_2 = t^{(m-1)} + \tau^{(m)}$ ,    $|\tau_2| \leq \text{eps} \cdot \text{ufp}(\tau_1)$ ,    $\tau_1 = \text{fl}(\tau_1 + \tau_2) = \text{fl}(t^{(m-1)} + \tau^{(m)}) = t^{(m)}$

*is satisfied for the final value of m. If $\sigma_{m-1} \leq \frac{1}{2}\text{eps}^{-1}\text{eta}$ is satisfied for the final value of m, then the vector $p^{(m)}$ is entirely zero. If $\sigma_{m-1} > \frac{1}{2}\text{eps}^{-1}\text{eta}$ is satisfied for the final value of m, then*

(3.5)                                    $\text{ufp}(\tau_1) \geq \Phi\sigma_{m-1}$ .

REMARK.    The assumption $\varrho \in \text{eps}\,\sigma_0\mathbb{Z}$ is necessary to ensure that $\text{ufp}(\varrho)$ is not too small: If Algorithm `Transform` does not stop for $m = 1$, then we will show that no rounding error occurs in the computation of $t^{(1)} = \text{fl}(t^{(0)} + \tau^{(1)})$. Without the assumption on $\varrho = t^{(0)}$ this need not be true if $0 < \varrho \leq \text{eps}|\tau^{(1)}|$.

PROOF OF LEMMA 3.4. For $\varrho = 0$, Algorithm `Transform` is identical to Algorithm 4.1 in Part I [22] of this paper with the parameter $\Phi = 2^{2M}\text{eps}$. Algorithm 4.1 in Part I is analyzed in Lemma 4.2 in Part I with the stronger assumption $2^{2M}\text{eps} \leq 1$.

Carefully going through the proof of Lemma 4.2 in Part I we identify the necessary changes to prove Lemma 3.4. We see that $2^M\text{eps} < 1$ suffices to guarantee that Algorithm 3.3 (`Transform`) stops. The assumption $\varrho \in \text{eps}\,\sigma_0\mathbb{Z}$ assures that (3.3) is satisfied for $m = 1$, and $t^{(0)} = \varrho$ gives (3.2) for $m = 1$. In the induction step we used $\sigma_{m-2} \geq \text{eps}^{-1}\text{eta}$ to verify that no rounding error has occurred in the previous computation of $2^{2M}\text{eps}\,\sigma_{m-2}$ in the "until"-condition. This is also true for $\Phi\sigma_{m-2}$ under our assumption $\text{eps} \leq \Phi \leq 1$. Next we proved

(3.6)                                    $|t^{(m-1)}| < \sigma_{m-2}$

in (4.6) in Part I to assure that no rounding error occurs in the computation of $t^{(m-1)} = \text{fl}(t^{(m-2)} + \tau^{(m-1)})$. But (3.6) follows by the previous "until"-condition on $|t^{(m-1)}|$ and $\text{eps} \leq \Phi \leq 1$. This proves (3.2) and (3.3) for all values of $m$, and (3.4) follows as well. If $\sigma_{m-1} \leq \frac{1}{2}\text{eps}^{-1}\text{eta}$ is satisfied for the final value of $m$, then (3.3) implies $|p_i^{(m)}| \leq \frac{1}{2}\text{eta}$, hence $p_i^{(m)} = 0$ for all $i$. Finally, $\text{eps} \leq \Phi \leq 1$ assures again that no rounding error occurred in the computation of $\Phi\sigma_{m-1}$ if $\sigma_{m-1} > \frac{1}{2}\text{eps}^{-1}\text{eta}$ and (3.5) follows. The lemma is proved. $\square$

A special application of Algorithm 3.3 (`Transform`) with $\varrho \neq 0$ might 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 $p_2, \cdots, p_n$ and not on $p_1$. We will use this in Section 6, where we will present Algorithm 6.4 (`AccSumK`) to compute a faithfully rounded result of $K$-fold accuracy.

The following lemma adapts some parts of the analysis of the original Algorithm 4.1 (`Transform`) in Part I and includes the offset $\varrho$. Note that the code in (3.7) without offset is the same as for Algorithm 4.5 (`AccSum`) in Part I.

LEMMA 3.5.    *Let p be a nonzero vector of n floating-point numbers and $\varrho \in \mathbb{F}$. Let* `res` *be computed as follows:*

(3.7)        $[\tau_1, \tau_2, p', \sigma] = \text{Transform}(p, \varrho)$        *% Parameter $\Phi$ replaced by $2^{2M}\text{eps}$*
$\qquad\qquad \text{res} = \text{fl}(\tau_1 + (\tau_2 + (\sum_{i=1}^{n} p_i')))$

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

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

$$(3.8) \qquad s + \varrho = \tau_1 + \tau_2 + \sum_{i=1}^{n} p_i' \quad \text{and} \quad \max |p_i'| \leq \mathtt{eps}\sigma \ ,$$

$$(3.9) \qquad \mathrm{fl}(\tau_1 + \tau_2) = \tau_1 \ , \quad \tau_1, \tau_2 \in \mathtt{eps}\sigma\mathbb{Z} \quad \text{and} \quad |\tau_2| \leq \mathtt{eps} \cdot \mathrm{ufp}(\tau_1) \ ,$$

$$(3.10) \qquad |s + \varrho - \mathtt{res}| < 2\mathtt{eps}(1 - 2^{-M-1})\mathrm{ufp}(\mathtt{res}) \ .$$

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

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

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

$$(3.11) \qquad \mathrm{ufp}(\tau_1) \geq 2^{2M}\mathtt{eps}\sigma \ .$$

*If the exponent* $2M$ *in the parameter* $\Phi$ *is changed into another integer, then* `res` *need not be a faithful rounding of* $s + \varrho$. *Abbreviate*

$$(3.12) \qquad \begin{array}{rcccl} \tau_3 & = & \mathrm{fl}(\sum_{i=1}^{n} p_i') & = & \sum_{i=1}^{n} p_i' - \delta_3 \ , \\ \tau_2' & = & \mathrm{fl}(\tau_2 + \tau_3) & = & \tau_2 + \tau_3 - \delta_2 \ , \\ \mathtt{res} & = & \mathrm{fl}(\tau_1 + \tau_2') & = & \tau_1 + \tau_2' - \delta_1 \ . \end{array}$$

*Then*

$$(3.13) \qquad s + \varrho = r + \delta \ \text{and} \ \mathtt{res} = \mathrm{fl}(r) \quad \text{for} \quad r := \tau_1 + \tau_2' \ \text{and} \ \delta := \delta_2 + \delta_3 \ .$$

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

$$(3.14) \qquad |\tau_3| = |\mathrm{fl}(\sum_{i=1}^{n} p_i')| \leq n\mathtt{eps}\sigma \ ,$$

$$(3.15) \qquad |\tau_2'| \leq (1 + \mathtt{eps})|\tau_2 + \tau_3| \leq (1 + \mathtt{eps})(\mathtt{eps} \cdot \mathrm{ufp}(\tau_1) + n\mathtt{eps}\sigma) < |\tau_1| \ ,$$

$$(3.16) \qquad |\mathtt{res}| \geq \frac{5}{8}|\tau_1| \ ,$$

$$(3.17) \qquad 2\mathtt{eps}^{-1}|\delta| < (1 - 2^{-M})|\mathtt{res}| < |\mathtt{res}| \ .$$

PROOF. The only difference to the assumptions of Lemma 4.3 in Part I of this paper is the additional parameter $\varrho$ in `Transform`. The proof and the assertions of Lemma 4.3 in Part I, however, are based on Lemma 4.2 in Part I and the error-free transformation $s = \tau_1 + \tau_2 + \sum p_i'$. With the offset $\varrho$ this changes into $s + \varrho = \tau_1 + \tau_2 + \sum p_i'$ as by Lemma 3.4. Carefully going through the proof of Lemma 4.3 in Part I we see that basically $s$ has to be replaced by $s + \varrho$, and the assertions (3.8) through (3.11) follow. The abbreviations (3.12) are the same as in (4.13) in the proof of Lemma 4.3 in Part I, implying (3.13). The remaining assertions (3.14) through (3.17) repeat corresponding assertions in the proof of Lemma 4.3 in Part I for the case $\sigma > \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta}$. $\qquad \square$

**4. The sign of a sum.** In Remark 2 following Lemma 4.3 in Part I of this paper we saw that the exponent in the lower bound $2^{2M}\mathtt{eps}\sigma_{m-1}$ in the "until"-condition cannot be decreased without jeopardizing faithful rounding. In the refined version of Algorithm 3.3 (`Transform`) we saw in Lemma 3.4 that most properties remain valid when replacing the parameter $\Phi$ by a factor as small as `eps`, i.e. the lower bound in the "until"-condition reads $\mathtt{eps}\sigma_{m-1}$.

The smaller $\Phi$, the less iterations are necessary. If only the sign of the sum $\sum p_i$ is to be computed and not necessarily a faithful rounding, then $\Phi$ may be decreased below $2^{2M}\mathtt{eps}$. Next we will show that if $\Phi$ is replaced by $2^M\mathtt{eps}$, then the sign of `res` as computed by Algorithm 4.5 (`AccSum`) in Part I is still equal to

the sign of the sum $\sum p_i$, and this not only for $n$ bounded by $2^{2M}\mathtt{eps} \leq 1$ but only requiring $2^M\mathtt{eps} < 1$. Moreover, $2^M\mathtt{eps}$ is the smallest possible choice for $\Phi$.

ALGORITHM 4.1. *Rigorous computation of* $\mathrm{sign}(\sum p_i)$, *also for huge lengths.*

$$\begin{aligned}
&\text{function } S = \mathtt{AccSign}(p)\\
&\quad [\tau_1, \tau_2, p', \sigma] = \mathtt{Transform}(p, 0) \qquad \text{\% \textit{Parameter} } \Phi \text{ \textit{replaced by} } 2^M\mathtt{eps}\\
&\quad S = \mathrm{sign}(\tau_1)
\end{aligned}$$

THEOREM 4.2. *Let $S$ be the result of Algorithm 4.1 (*$\mathtt{AccSign}$*) applied to a vector of floating-point numbers $p_i, 1 \leq i \leq n$. Assume that $2^M\mathtt{eps} < 1$ is satisfied for $M := \lceil \log_2(n+2) \rceil$.*

*Then Algorithm $\mathtt{AccSign}$ will stop, and*

$$(4.1) \qquad\qquad S = \mathrm{sign}\Big(\sum_{i=1}^{n} p_i\Big) \; .$$

*The parameter $\Phi$ cannot be replaced by a power of $2$ smaller than $2^M\mathtt{eps}$ without jeopardizing (4.1). The algorithm needs $(4m+2)n + \mathcal{O}(m)$ flops for $m$ "repeat-until"-loops in $\mathtt{Transform}$.*

PROOF. Define $s := \sum_{i=1}^{n} p_i$. Without loss of generality assume the vector $p$ to be nonzero. The assumptions of Lemma 3.4 are satisfied, and (3.2), (3.4) and $\varrho = 0$ imply

$$(4.2) \qquad\qquad s = \tau_1 + \tau_2 + \sum_{i=1}^{n} p'_i \; .$$

If $\sigma \leq \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta}$, then $p'$ is entirely zero and $s = \tau_1 + \tau_2$, so that $\tau_1 = \mathrm{fl}(\tau_1 + \tau_2) = \mathrm{fl}(s)$.

If $\sigma > \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta}$, then $\tau_1 = t^{(m)}$ and the "until"-condition in Algorithm 3.3 ($\mathtt{Transform}$) imply $|\tau_1| \geq 2^M\mathtt{eps}\sigma$, so that $|\tau_2| \leq \mathtt{eps}|\tau_1|$ and (3.3) give

$$|\tau_2| + \Big|\sum_{i=1}^{n} p'_i\Big| \leq \mathtt{eps}|\tau_1| + (2^M - 2)\mathtt{eps}\sigma < (1 - \mathtt{eps})|\tau_1|$$

and $\mathrm{sign}(s) = \mathrm{sign}(\tau_1)$ by (4.2).

To see that $\Phi$ is optimal, consider

$$(4.3) \qquad p^{(0)} = [c \quad c \quad c \quad c \quad -\tfrac{1}{8}\mathtt{eps}^{-1} \quad -4 + 4\mathtt{eps}] \in \mathbb{F}^6 \quad \text{ with } \quad c := \frac{1}{32}\mathtt{eps}^{-1} + 1$$

in a floating-point format with relative rounding error unit $\mathtt{eps} \leq \frac{1}{64}$. Then $M = 3, \mu = \frac{1}{8}\mathtt{eps}^{-1}$ and $\sigma = \mathtt{eps}^{-1}$. Hence rounding tie to even implies $p^{(1)} = [1 \quad 1 \quad 1 \quad 1 \quad 0 \quad 4\mathtt{eps}]$ and $\tau^{(1)} = 4 \cdot \frac{1}{32}\mathtt{eps}^{-1} - \frac{1}{8}\mathtt{eps}^{-1} - 4 = -4$ in the notation of Algorithm 3.3 ($\mathtt{Transform}$). With the parameter $\Phi = 2^{M-1}\mathtt{eps}$ the first inequality in the "until"-condition would read $|t^{(m)}| \geq 4$, so the "repeat-until"-loop would be finished. However, $t^{(0)} = 0$, so $\tau_1 = -4$ but $s = +4\mathtt{eps}$. $\qquad\square$

Note that for the example in (4.3) with the parameter $\Phi = 2^{M-1}\mathtt{eps}$ in $\mathtt{AccSign}$ it would also not help to compute $\mathtt{res}$ as in Algorithm 4.5 ($\mathtt{AccSum}$) in Part I because $\tau_2 = 0$ and $\mathrm{fl}(\sum p_i^{(1)}) = 4$ by rounding tie to even, so that $\mathtt{res} = 0$ but $s = 4\mathtt{eps}$.

The improved version for sign determination is applicable in single precision for dimensions up to $8.3 \cdot 10^6$, and in double precision for dimensions up to $4.5 \cdot 10^{15}$.

**5. $K$-fold faithful rounding.** 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 in [2, 3, 7, 6]. 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 [14]. For two and more pieces this technique was used in [21] to compute highly accurate inclusions of the solution of systems of linear equations. Later this technique was also called *staggered correction*.

To gain as much accuracy as possible, the pieces may be required to be non-overlapping, i.e. their bit representations should have no bit in common [19, 20]. Another approach is to assume that two numbers may only overlap in bits where those of the first one are zero [24].

In this section we define *$K$-fold faithful rounding*. Recall that by Definition 2.2 a floating-point number $f \in \mathbb{F}$ is called a faithful rounding of a real number $r \in \mathbb{R}$ if

$$(5.1) \qquad \mathrm{pred}(f) < r < \mathrm{succ}(f) \ .$$

Before extending this definition to $K$-fold faithful rounding, we collect some properties of the ordinary faithful rounding.

LEMMA 5.1. *Let $r \in \mathbb{R}$ be given, and let $f \in \mathbb{F}$ be a faithful rounding of $r$. Then*

$$(5.2) \qquad f \notin \mathbb{U} \quad \Rightarrow \quad |r - f| < 2\mathtt{eps} \cdot \mathrm{ufp}(r) \quad and \quad |r - f| < 2\mathtt{eps} \cdot \mathrm{ufp}(f)$$

*and*

$$(5.3) \qquad f \in \mathbb{U} \quad \Rightarrow \quad |r - f| < \mathtt{eta} \ .$$

PROOF. If $f \in \mathbb{U}$, then $f \pm \mathtt{eta}$ are the neighbors of $f$, and the definition of faithful rounding (5.1) yields $|r - f| < \mathtt{eta}$. This proves (5.3).

To prove (5.2) we assume $f \notin \mathbb{U}$ and without loss of generality $f > 0$. For $f$ not being a power of 2 we have $\mathrm{ufp}(f) \le \mathrm{ufp}(r)$, so (5.1) and (2.11) in Lemma 2.1 imply

$$|r - f| < 2\mathtt{eps} \cdot \mathrm{ufp}(f) \le 2\mathtt{eps} \cdot \mathrm{ufp}(r) \ .$$

If $f$ is a power of 2, then this is also true if $f \le r < \mathrm{succ}(f)$ because $\mathrm{succ}(f) = (1 + 2\mathtt{eps})\mathrm{ufp}(f)$. If $(1 - \mathtt{eps})f = \mathrm{pred}(f) < r < f$ then

$$|r - f| = f - r < \mathtt{eps} \cdot f = \mathtt{eps} \cdot \mathrm{ufp}(f) = 2\mathtt{eps} \cdot \mathrm{ufp}(r) \ .$$

The lemma is proved.                                                                                 □

Lemma 5.1 is formulated for general $r \in \mathbb{R}$. For our main application, the approximation of the sum $s = \sum p_i$ of floating-point numbers $p_i$, we have $p_i \in \mathtt{eta}\mathbb{Z}$ and therefore $s \in \mathtt{eta}\mathbb{Z}$. Thus (5.3) reads in this case $f \in \mathbb{U} \Rightarrow r = f$. Before we define $K$-fold faithful rounding, we introduce the concept of a non-overlapping sequence of floating-point numbers.

DEFINITION 5.2. *A sequence of floating-point numbers $f_1, \cdots, f_k \in \mathbb{F}$ is called* non-overlapping*, if*

$$(5.4) \qquad \begin{aligned} |f_{i+1}| &< 2\mathtt{eps} \cdot \mathrm{ufp}(f_i) \qquad && if \quad f_i \ne 0 \\ f_{i+1} &= 0 && if \quad f_i = 0 \end{aligned}$$

*holds true for $1 \le i < k$.*

REMARK. Note that trailing zeros are possible, and that a sequence of zeros is by definition non-overlapping.

The definition implies that for $f_i f_{i+1} \neq 0$ the binary expansions of $f_i$ and $f_{i+1}$ do not have a bit in common. The following properties hold true.

LEMMA 5.3. *Let $f_1, \cdots, f_k$ be a non-overlapping sequence of floating-point numbers. Then*

$$(5.5) \qquad\qquad \mathrm{ufp}(f_k) \leq \mathsf{eps}^{k-1}\mathrm{ufp}(f_1) \ .$$

*Moreover,*

$$(5.6) \qquad\qquad |\sum_{\nu=1}^{k} f_\nu| < 2\mathrm{ufp}(f_1) \ .$$

REMARK. Note that Lemma 5.3 is applicable to any subsequence $f_{i_1}, \cdots, f_{i_m}$ with $i_1 > \cdots > i_m$ because it is non-overlapping as well.

PROOF. Let $(f, g)$ be a non-overlapping sequence. Then the definition (5.4) of non-overlapping implies

$$(5.7) \qquad\qquad \mathrm{ufp}(g) \leq \mathsf{eps} \cdot \mathrm{ufp}(f) \ ,$$

and (5.5) follows. To prove (5.6) we first show

$$(5.8) \qquad\qquad |f| + 2\mathrm{ufp}(g) \leq 2\mathrm{ufp}(f) \ .$$

If $f \in \mathbb{U}$, then the definition of the set $\mathbb{U}$ and (5.5) imply $\mathrm{ufp}(g) \leq \frac{1}{2}\mathsf{eta}$ and therefore $g = 0$. If $f \notin \mathbb{U}$, then $|f| \leq \mathrm{pred}(2\mathrm{ufp}(f)) = 2(1 - \mathsf{eps})\mathrm{ufp}(f)$, and (5.7) implies (5.8). Without loss of generality we can omit trailing zeros in the sequence $f_1, \cdots, f_k$ and assume $f_k \neq 0$. Then

$$|\sum_{\nu=1}^{k} f_\nu| < \sum_{\nu=1}^{k-1} |f_\nu| + 2\mathrm{ufp}(f_k) \leq \sum_{\nu=1}^{k-2} |f_\nu| + 2\mathrm{ufp}(f_{k-1}) \leq \cdots \leq 2\mathrm{ufp}(f_1) \ ,$$

proving (5.6). $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

Our aim is a sequence $(f_1, \cdots, f_k)$ such that $\sum f_\nu$ has a small relative error with respect to $s$, i.e.

$$(5.9) \qquad\qquad |s - \sum_{\nu=1}^{k} f_\nu| < 2\mathsf{eps}^k|s| \ .$$

For that we extend the Definition 2.2 of ordinary faithful rounding to a sequence $f_1, \cdots, f_k$.

DEFINITION 5.4. *A sequence $f_1, \cdots, f_k \in \mathbb{F}$ is called a (k-fold) faithful rounding of $s \in \mathbb{R}$ if*

$$(5.10) \qquad\qquad f_i \in \square(s - \sum_{\nu=1}^{i-1} f_\nu) \qquad \text{for } 1 \leq i \leq k \ .$$

*The sequence is called a strongly faithful rounding if it is in addition non-overlapping.*

Of course, a strongly faithful rounding is faithful, but the converse is not true. However, the following lemma shows that the members of a sequence representing a faithful rounding may overlap in at most one bit, and subnormal numbers can only occur at the end.

LEMMA 5.5. *Suppose the sequence $f_1, \cdots, f_k \in \mathbb{F}$ is a faithful rounding of some $s \in \mathbb{R}$. If $|f_{i-1}| \geq \frac{1}{2}\mathsf{eps}^{-1}\mathsf{eta}$ for $1 < i \leq k$, then*

$$(5.11) \qquad |f_i| \leq 2\mathsf{eps} \cdot \mathrm{ufp}(f_{i-1}) \qquad \text{with equality if and only if} \quad \mathrm{ufp}(f_i) = |f_i| \ .$$

*Moreover, $|f_m| > 2\mathsf{eta}$ implies $f_i \notin \mathbb{U}$ for $1 \leq i < m$, and $|f_m| < \frac{1}{2}\mathsf{eps}^{-1}\mathsf{eta}$ implies $f_i = 0$ for $m < i \leq k$.*

PROOF. Let $f_1, \cdots, f_k \in \mathbb{F}$ be a faithful rounding of $s \in \mathbb{R}$. Then by definition

$$(5.12) \qquad \sum_{\nu=1}^{m-1} f_\nu + \mathrm{pred}(f_m) < s < \sum_{\nu=1}^{m-1} f_\nu + \mathrm{succ}(f_m) \qquad \text{for } 1 \leq m \leq k .$$

Assume $f_i \geq 0$. Then using (5.12) for $m = i$ and $m = i - 1$ gives

$$\mathrm{pred}(f_i) < s - \sum_{\nu=1}^{i-1} f_\nu < \mathrm{succ}(f_{i-1}) - f_{i-1} \in \mathbb{F} ,$$

so

$$\mathrm{ufp}(f_i) \leq f_i \leq \mathrm{succ}(f_{i-1}) - f_{i-1} \leq 2\mathsf{eps} \cdot \mathrm{ufp}(f_{i-1}) ,$$

where the last inequality follows by the assumption $|f_{i-1}| \geq \frac{1}{2}\mathsf{eps}^{-1}\mathsf{eta}$. The case $f_i < 0$ is treated similarly, and this proves (5.11), and also the very last statement of the lemma. For $f_m \geq 0$ and $f_m > 2\mathsf{eta}$, (5.12) yields

$$\mathrm{succ}(f_{m-1}) - f_{m-1} > s - \sum_{\nu=1}^{m-2} f_\nu - f_{m-1} > \mathrm{pred}(f_m) \geq 2\mathsf{eta} .$$

Hence $f_{m-1} \notin \mathbb{U}$, and an induction argument finishes the proof. $\qquad \square$

Next we show that faithful rounding is sufficient to satisfy our anticipated bound (5.9) on the relative error with respect to $s$.

PROPOSITION 5.6. *If a sequence $f_1, \cdots, f_k \in \mathbb{F}$ is a faithful rounding of some $s \in \mathbb{R}$, then*

$$(5.13) \qquad \left| s - \sum_{\nu=1}^{k} f_\nu \right| \quad < \quad \begin{cases} 2\mathsf{eps}^k \mathrm{ufp}(s) & \text{if } f_k \notin \mathbb{U} \\ 2\mathsf{eps}^k \mathrm{ufp}(f_1) & \text{if } f_k \notin \mathbb{U} \\ \mathsf{eta} & \text{if } f_k \in \mathbb{U} . \end{cases}$$

*Moreover, a sequence representing a faithful rounding need not to be strongly faithful, i.e. may be overlapping.*

PROOF. For $0 \leq i \leq k$ define $\Delta_i := s - \sum_{\nu=1}^{i} f_\nu$. By assumption, $f_i \in \square(\Delta_{i-1})$ for $1 \leq i \leq k$. If $f_k \in \mathbb{U}$, then (5.13) follows by Lemma 5.1. Suppose $f_k \notin \mathbb{U}$, so that Lemma 5.5 implies $f_i \notin \mathbb{U}$ for $1 \leq i \leq k$. Now Lemma 5.1 yields $|\Delta_i| < 2\mathsf{eps} \cdot \min\bigl(\mathrm{ufp}(\Delta_{i-1}), \mathrm{ufp}(f_i)\bigr)$, hence

$$(5.14) \qquad \mathrm{ufp}(\Delta_i) \leq \mathsf{eps} \cdot \min\bigl(\mathrm{ufp}(\Delta_{i-1}), \mathrm{ufp}(f_i)\bigr)$$

and

$$\mathrm{ufp}(\Delta_k) \leq \mathsf{eps}^{k-1}\mathrm{ufp}(\Delta_1) < 2\mathsf{eps}^k \cdot \min\bigl(\mathrm{ufp}(s), \mathrm{ufp}(f_1)\bigr) .$$

This proves (5.13) because the rightmost expression is a power of 2. Finally consider

$$f_1 = 1 - \mathsf{eps}, \quad f_2 = \mathsf{eps} \quad \text{and} \quad s = 1 - \mathsf{eps}^3 ,$$

then

$$f_1 + \mathrm{pred}(f_2) = 1 - \mathsf{eps}^2 \quad < \quad s \quad < 1 + 2\mathsf{eps}^2 = f_1 + \mathrm{succ}(f_2)$$
$$\mathrm{pred}(f_1) = 1 - 2\mathsf{eps} \quad < \quad s \quad < 1 = \mathrm{succ}(f_1)$$

implies that $(f_1, f_2)$ is a faithful rounding of $s$. However, it is overlapping by $|f_2| \geq 2\mathsf{eps} \cdot \mathrm{ufp}(f_1) = \mathsf{eps}$. $\square$

Note the difference to Algorithm 4.8 (SumK) in [18]. There, the result can also be stored in an array of length $K$ and it is of a quality *as if* computed in $K$-fold precision, i.e. with relative rounding error unit $\mathsf{eps}^K$. So the *precision* is $K$-fold, the *accuracy*, however, depends on the condition number of the sum. The accuracy of the $K$-fold faithfully rounded result is independent of the condition number.

It is straightforward to derive from Algorithm 4.5 (AccSum) in Part I a new algorithm computing a result *as if* computed in $K$-fold precision, similar to SumK in [18]. This fact has been mentioned independently to the first author by Nicolas Louvet [15].

**6. Algorithms for $K$-fold accuracy.** Let floating-point numbers $p_1, \cdots, p_n$ be given. Next we will derive an algorithm to compute a sequence $f_1, \cdots, f_K \in \mathbb{F}$ representing a strongly faithful rounding of the sum $s := \sum p_i$.

LEMMA 6.1. *Let for a nonzero vector $p$ of $n$ floating-point numbers and $\varrho \in \mathbb{F}$ the assumptions of Lemma 3.5 be satisfied. Assume* res $\in \mathbb{F}$ *has been computed with the code given in (3.7) in Lemma 3.5. Set $f_1 := $ res, so that $f_1 \in \square(s + \varrho)$. Furthermore, suppose $f_2 \in \square(s + \varrho - f_1)$ for some $f_2 \in \mathbb{F}$. Then $f_1, f_2$ are non-overlapping.*

PROOF. If $f_2 \notin \mathbb{U}$, then (5.2) and (3.10) imply

$$
\begin{aligned}
|f_2| &\leq |s + \varrho - f_1| + |s + \varrho - f_1 - f_2| \leq |s + \varrho - f_1| + 2\mathsf{eps} \cdot \mathrm{ufp}(s + \varrho - f_1) \\
&\leq (1 + 2\mathsf{eps})|s + \varrho - f_1| < 2\mathsf{eps}(1 + 2\mathsf{eps})(1 - 2^{-M-1})\mathrm{ufp}(f_1) \\
&< 2\mathsf{eps} \cdot \mathrm{ufp}(f_1) \ ,
\end{aligned}
$$

proving that $f_1$ and $f_2$ are non-overlapping. If $f_2 \in \mathbb{U}$, then $f_2 \in \square(s + \varrho - f_1)$ implies $f_2 = s + \varrho - f_1$, and the assertion follows similarly. $\qquad\square$

When applying the code given in Lemma 3.5 to a vector $p$ and $\varrho \in \mathbb{F}$, it calculates a faithful rounding res of $s + \varrho$, and also extracts $p$ into some vector $p'$. However, we have no *equation* relating $s$, $\varrho$, res and $p'$, but only an *estimation* which follows from res $\in \square(s + \varrho)$. The only *equation* between the mentioned quantities we know up to now is $s + \varrho = \tau_1 + \tau_2 + \sum p_i'$ as by Lemma 3.5.

However, to be able to repeatedly apply Lemma 6.1 to produce a non-overlapping sequence $(f_1, \cdots, f_K)$, we need a faithful rounding of $s + \varrho - $ res. This can be done by applying the code given in (3.7) in Lemma 3.5 to $\varrho$ and the vector $p$ appended by $-$res. However, this would extract the entire vector $[p, -$res$]$ again and is ineffective. In the following Algorithm 6.2 (TransformK) we give an effective way to compute a single floating-point number $R$ relating $s$, $\varrho$, res and $p'$ by an equation. It paves the way to compute a strongly faithful sequence of $K$ numbers $f_1, \cdots, f_K$ approximating $s$, thus establishing an approximation of $K$-fold accuracy.

ALGORITHM 6.2. *Error-free vector transformation including faithful rounding.*

$$
\begin{aligned}
&\text{function } [\mathtt{res}, R, p'] = \mathtt{TransformK}(p, \varrho) \\
&\quad [\tau_1, \tau_2, p', \sigma] = \mathtt{Transform}(p, \varrho) \qquad\quad \text{\% Parameter } \Phi \text{ replaced by } 2^{2M}\mathsf{eps} \\
&\quad \mathtt{res} = \mathrm{fl}(\tau_1 + (\tau_2 + (\textstyle\sum_{i=1}^n p_i'))) \\
&\quad R = \mathrm{fl}(\tau_2 - (\mathtt{res} - \tau_1))
\end{aligned}
$$

LEMMA 6.3. *Let $p$ be a nonzero vector of $n$ floating-point numbers and let $\varrho \in \mathbb{F}$. Abbreviate $s := \sum p_i$. Define $M := \lceil \log_2(n+2) \rceil$, and assume that $2^{2M}\mathsf{eps} \leq 1$. Furthermore, assume that $\varrho \in \mathsf{eps}\sigma_0\mathbb{Z}$ is satisfied for $\mu := \max_i |p_i|$ and $\sigma_0 = 2^{M + \lceil \log_2 \mu \rceil}$. Let* res, $R$ *and $p'$ be the results of Algorithm 6.2 (TransformK) applied to $p$ and $\varrho$. Then* res *is a faithful rounding of $s + \varrho$ and*

$$
\tag{6.1} s + \varrho - \mathtt{res} = R + \sum_{i=1}^n p_i' \ .
$$

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

$$
\tag{6.2} R \in \mathsf{eps}\sigma'\mathbb{Z}
$$

*is satisfied for $\mu' := \max_i |p_i'|$ and $\sigma' := 2^{M + \lceil \log_2 \mu' \rceil}$. Moreover, no rounding error occurs in the computation of $R$, i.e.*

$$
\tag{6.3} \Delta := \mathtt{res} - \tau_1 = \mathrm{fl}(\mathtt{res} - \tau_1) \quad and \quad R = \tau_2 - \Delta = \mathrm{fl}(\tau_2 - \Delta) \ .
$$

*Algorithm 6.2 (*`TransformK`*) requires $(4m+3)n+\mathcal{O}(m)$ flops for $m$ executions of the "repeat-until"-loop in* `Transform`*.*

REMARK 1. The main purpose of Algorithm `TransformK` is to replace the pair $\tau_1, \tau_2$ by the pair `res`, $R$ without error. This makes it suitable for cascading, see Algorithm 6.4 (`AccSumK`).

REMARK 2. The main part of the proof of Lemma 6.3 will be to show that no rounding error can occur in the computation of $R$, that is $R = \tau_2 - (\texttt{res} - \tau_1)$. The proof is involved and moved to the Appendix.

REMARK 3. As is Matlab convention, an output parameter, in this case $M$ for `Transform`, has been omitted. The output parameter $\sigma$ is not needed but added for clarity.

PROOF OF LEMMA 6.3. For the computation of $\tau_1, \tau_2, p', \sigma$ and `res`, Algorithm 6.2 (`TransformK`) uses the same piece of code as in (3.7) in Lemma 3.5, so the assertions of Lemma 3.5 are valid. Hence the assertions follow if $\texttt{res} = 0$.

First assume $\sigma \leq \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$. Then $p_i' = 0$ for all $i$ by Lemma 3.5, and (3.8) and (3.9) imply $s + \varrho = \tau_1 + \tau_2$ and $\texttt{res} = \text{fl}(\tau_1 + \tau_2) = \tau_1$, so $R = \tau_2$ and (6.1) and (6.3) follow. The case $\sigma \leq \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$ is finished.

Henceforth we assume $\sigma > \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$. We use the notation in Lemma 3.5, especially (3.12). Then $|\tau_2'| < |\tau_1|$ by (3.15), and $\texttt{res} = \text{fl}(\tau_1 + \tau_2')$ and (2.14) in Lemma 2.5 imply

$$(6.4) \qquad \Delta := \texttt{res} - \tau_1 = \text{fl}(\texttt{res} - \tau_1) \in \mathbb{F} \ .$$

Next we have to show $R \in \mathbb{F}$, i.e.

$$(6.5) \qquad R = \text{fl}(\tau_2 - \Delta) = \tau_2 - \Delta$$

provided $\sigma > \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$. This proof is quite involved and left to the Appendix. It verifies (6.3). Now (6.4) and (6.5) yield $R = \tau_1 + \tau_2 - \texttt{res}$, and (3.8) gives

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

This proves (6.1). To see (6.2) note that (2.8), (3.11) and (2.4) imply $\Delta = \texttt{res} - \tau_1 \in \texttt{eps} \cdot \text{ufp}(\tau_1)\mathbb{Z} \subseteq 2^{2M}\texttt{eps}^2\sigma\mathbb{Z}$. Furthermore, (3.9) and $2^{2M}\texttt{eps} \leq 1$ yield $\tau_2 \in \texttt{eps}\sigma\mathbb{Z} \subseteq 2^{2M}\texttt{eps}^2\sigma\mathbb{Z}$, and (2.6) gives

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

But (3.8) and the definition of $\sigma'$ imply $\texttt{eps}\sigma' \leq 2^M\texttt{eps}^2\sigma$, and (6.2) follows by (2.4). $\qquad \square$

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

ALGORITHM 6.4. *Accurate summation computing $K$-fold faithful rounding.*

$$
\begin{aligned}
&\text{function } \texttt{Res} = \texttt{AccSumK}(p, K) \\
&\quad p^{(0)} = p, \ R_0 = 0 \\
&\quad \text{for } k = 1 : K \\
&\qquad [\texttt{Res}_k, R_k, p^{(k)}] = \texttt{TransformK}(p^{(k-1)}, R_{k-1}) \\
&\qquad \text{if } \texttt{Res}_k \in \mathbb{U}, \ \texttt{Res}_{k+1..K} = 0, \text{ return; end if} \\
&\quad \text{end for}
\end{aligned}
$$

PROPOSITION 6.5. *Let $p$ be a vector of $n$ floating-point numbers and abbreviate $s := \sum p_i$. Define $M := \lceil \log_2(n+2) \rceil$, and assume that $2^{2M}\texttt{eps} \leq 1$. Let $1 \leq K \in \mathbb{N}$ be given, and let `Res` be the result vector of Algorithm 6.4 (`AccSumK`) applied to $p$ and $K$.*

*Then* $\mathtt{Res}_1, \cdots, \mathtt{Res}_K$ *is a strongly faithful rounding of* $s$. *Furthermore,*

$$(6.6) \qquad s = \sum_{\nu=1}^{k} \mathtt{Res}_\nu + R_k + \sum_{i=1}^{n} p_i^{(k)} \quad for\ 1 \le k \le K .$$

*Abbreviate* $\overline{\mathtt{Res}} := \sum_{k=1}^{K} \mathtt{Res}_k$. *Then* $\mathtt{Res}_k \in \mathbb{U}$ *for some* $1 \le k \le K$ *implies* $\overline{\mathtt{Res}} = s$, *and* $\mathtt{Res}_k \notin \mathbb{U}$ *implies*

$$(6.7) \qquad |s - \overline{\mathtt{Res}}| < 2\mathtt{eps}^K \mathrm{ufp}(s) \le 2\mathtt{eps}^K |s|$$

*and*

$$(6.8) \qquad |s - \overline{\mathtt{Res}}| < 2\mathtt{eps}^K \mathrm{ufp}(\mathtt{Res}_1) \le 2\mathtt{eps}^K |\mathtt{Res}_1| .$$

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

REMARK 2. For $k \ge 3$, say, it might be advantageous to eliminate zero summands in the vectors $p^{(k)}$. Assuming that it seems unlikely that extracted vectors $p^{(k)}$ are ill-conditioned for $k \ge 1$ and using the fact that the new $\sigma$ needs not to be computed in subsequent calls to $\mathtt{TransformK}$, $\mathtt{AccSumK}$ needs $(4m + 5K + 3)n + \mathcal{O}(m + K)$ flops if the "repeat-until"-loop in the first extraction $\mathtt{TransformK}$ is executed $m$ times.

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 $s \in \mathtt{eta}\mathbb{Z}$, and also $s - \sum f_k \in \mathtt{eta}\mathbb{Z}$ for arbitrary floating-point numbers $f_k$. So for a certain $K$, the sum $s$ is *always* stored exactly in $\sum_{k=1}^{K} \mathtt{Res}_k$. The maximum $K$ can be estimated solely by the exponent range because the sequence $(\mathtt{Res}_1, \cdots, \mathtt{Res}_K)$ is non-overlapping. In IEEE 754 single precision at most $K \le 12$ summands are needed, in IEEE 754 double precision at most $K \le 40$.

REMARK 4. In an actual implementation, of course, $R$ and the vector $p$ can be reused. So, for example, the main statement in Algorithm 6.4 ($\mathtt{AccSumK}$) reads $[\mathtt{Res}_k, R, p] = \mathtt{TransformK}(p, R)$.

PROOF OF PROPOSITION 6.5. For zero input vector $p$ the assertions are evident; henceforth we assume $p$ to be nonzero.

For $k = 0$ the assumptions of Lemma 6.3 are satisfied, so $\mathtt{Res}_1$ is a faithful rounding of $s = \sum p_i^{(0)}$. Moreover, by (6.2) the assumptions of Lemma 6.3 are satisfied for $k \ge 1$ as well. With $s_k := \sum p_i^{(k)}$ this means

$$
\begin{aligned}
s \;=\; s_0 \;&=\; \mathtt{Res}_1 + R_1 + s_1 \\
&=\; \mathtt{Res}_1 + \mathtt{Res}_2 + R_2 + s_2 \\
&=\; \cdots \\
&=\; \textstyle\sum_{k=1}^{K} \mathtt{Res}_k + R_K + s_K ,
\end{aligned}
$$

which is true because $\mathtt{Res}_k \in \mathbb{U}$ implies $R_k$ and all $p_i^{(k)}$ to be zero by Lemma 6.3. This proves (6.6), and that $\mathtt{Res}_k$ is a faithful rounding of $s - \sum_{\nu=1}^{k-1} \mathtt{Res}_\nu$ for $1 \le k \le K$. Thus the sequence $\mathtt{Res}_1, \cdots, \mathtt{Res}_K$ is a faithful rounding of $s$, and Proposition 5.6 proves (6.7) and (6.8). Moreover, Lemma 6.1 implies $\mathtt{Res}_1, \cdots, \mathtt{Res}_K$ to be non-overlapping, thus $\mathtt{Res}_1, \cdots, \mathtt{Res}_K$ is even a strongly faithful rounding of $s$. $\qquad\square$

Now it becomes clear why so much effort was spent (see the Appendix) to prove (6.1) and (6.2): It is the key to cascade $\mathtt{TransformK}$ in Algorithm 6.4 ($\mathtt{AccSumK}$) without increasing the length of the input vector. Especially for large gaps in the input vector, $R_{k-1}$ in $\mathtt{AccSumK}$ may be large compared to $p^{(k-1)}$. In that case a number of case distinctions were necessary in the proof in the Appendix to assure that no rounding error can occur in the computation of $R_k$ in $\mathtt{TransformK}$.

**7. Rounding to nearest and directed rounding.** Using the results of the previous section we can derive algorithms for computing $s := \sum p_i$ with directed rounding, i.e. $\texttt{resD} := \max\{f \in \mathbb{F} : f \leq s\}$ for rounding downwards, and $\texttt{resU} := \min\{f \in \mathbb{F} : s \leq f\}$ for rounding upwards. Note that $\texttt{resD} = \texttt{resU}$ is equivalent to $s \in \mathbb{F}$. As an example we display Algorithm 7.1 ($\texttt{DownSum}$) for rounding downwards, its counterpart looks similarly. The proof of correctness follows straightforwardly by Lemma 6.3 and the definition of faithful rounding.

ALGORITHM 7.1. *Accurate summation with rounding downwards.*

$$
\begin{aligned}
&\text{function } \texttt{resD} = \texttt{DownSum}(p) \\
&\quad [\texttt{res}, R, p'] = \texttt{TransformK}(p, 0) \quad \% \; s - \texttt{res} = R + \sum p_i' \\
&\quad \delta = \texttt{TransformK}(p', R) \qquad\quad \% \; \delta \in \square(s - \texttt{res}) \\
&\quad \text{if } \delta < 0, \; \texttt{resD} = \text{pred}(\texttt{res}) \quad \% \; s < \texttt{res} \\
&\quad \text{else } \texttt{resD} = \texttt{res} \qquad\qquad\quad \% \; \texttt{res} \leq s \\
&\quad \text{end if}
\end{aligned}
$$

Next we show how to compute the rounded to nearest result of the sum of a vector of floating-point numbers. If $\texttt{res}$ is a faithful rounding of $s := \sum p_i$ and $\texttt{res} \neq s$, then $\texttt{res}$ must be one of the immediate floating-point neighbors of $s$. Define for $c \in \mathbb{F} \setminus \{\pm\infty\}$ with finite predecessor and successor

$$
(7.1) \qquad \mathcal{M}^-(c) := \frac{1}{2}\big(\text{pred}(c) + c\big) \quad \text{and} \quad \mathcal{M}^+(c) := \frac{1}{2}\big(c + \text{succ}(c)\big).
$$

Then $\mathcal{M}^-(\texttt{res}) < s < \mathcal{M}^+(\texttt{res})$ implies $\text{fl}(s) = \texttt{res}$,

$$
(7.2) \qquad
\begin{aligned}
&s < \mathcal{M}^-(\texttt{res}) \Rightarrow \text{fl}(s) = \text{pred}(\texttt{res}) \quad \text{and} \quad \mathcal{M}^+(\texttt{res}) < s \Rightarrow \text{fl}(s) = \text{succ}(\texttt{res}) \;, \\
&s < \mathcal{M}^+(\texttt{res}) \Rightarrow \text{fl}(s) \neq \text{succ}(\texttt{res}) \quad \text{and} \quad \mathcal{M}^-(\texttt{res}) < s \Rightarrow \text{fl}(s) \neq \text{pred}(\texttt{res}) \;,
\end{aligned}
$$

and for the cases $\mathcal{M}^-(\texttt{res}) = s$ and $s = \mathcal{M}^+(\texttt{res})$ the rounding depends on the tie. So if $\text{sign}(\texttt{res} - s)$ is known, then the rounding to nearest $\text{fl}(s)$ of $s$ can be calculated by computing a faithful rounding of $s - \mu$ with $\mu \in \{\mathcal{M}^-(\texttt{res}), \mathcal{M}^+(\texttt{res})\}$, and it can be decided which of $\mathcal{M}^-$ or $\mathcal{M}^+$ to choose for $\mu$. This is true because a faithful rounding of $s - \mu$ detects the sign with certainty. The sign of $\texttt{res} - s$ could be calculated by calling $\texttt{AccSum}(s, -\texttt{res})$ using Algorithm 4.5 ($\texttt{AccSum}$) in Part I. However, we can do much better. The following lemma shows how enough information can be extracted from $\texttt{AccSum}$ to decide which of $\text{pred}(\texttt{res})$ or $\text{succ}(\texttt{res})$ the true rounding to nearest $\text{fl}(s)$ can only be in case $\text{fl}(s) \neq \texttt{res}$.

LEMMA 7.2. *Let $p$ be a vector of $n$ floating-point numbers, define $M := \lceil \log_2(n+2) \rceil$ and assume $2^{2M}\texttt{eps} \leq 1$. Let $\texttt{res}$ and $\delta_1$ be computed as follows:*

$$
(7.3) \qquad
\begin{aligned}
&[\tau_1, \tau_2, p'] = \texttt{Transform}(p, 0) \\
&\tau_2' = \text{fl}(\tau_2 + (\textstyle\sum_{i=1}^n p_i')) \qquad \% \; \textit{Parameter } \Phi \textit{ replaced by } 2^{2M}\texttt{eps} \\
&[\texttt{res}, \delta_1] = \texttt{FastTwoSum}(\tau_1, \tau_2')
\end{aligned}
$$

*Abbreviate $s := \sum_{i=1}^n p_i$. Then $\delta_1 = 0$ implies $\text{fl}(s) = \texttt{res}$, and*

$$
(7.4) \qquad \text{fl}(s) \in \left\{
\begin{array}{ll}
\{\text{pred}(\texttt{res}), \texttt{res}\} & \text{if } \delta_1 < 0 \\
\{\texttt{res}, \text{succ}(\texttt{res})\} & \text{if } \delta_1 > 0 \;.
\end{array}
\right.
$$

PROOF. The first statement in Algorithm 2.4 ($\texttt{FastTwoSum}$) is $\texttt{res} = \text{fl}(\tau_1 + \tau_2')$, so the result $\texttt{res}$ computed in (7.3) is the same as in (3.7) for $\varrho = 0$. Moreover, the assumptions of Lemma 3.5 are satisfied, so $\texttt{res}$ is a faithful rounding of $s$. This means $\text{fl}(s) \in \{\text{pred}(\texttt{res}), \texttt{res}, \text{succ}(\texttt{res})\}$. If the final value of $\sigma$ in $\texttt{Transform}$ satisfies $\sigma \leq \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$, then the vector $p'$ is entirely zero and $\text{fl}(s) = \text{fl}(\tau_1 + \tau_2) = \texttt{res}$. Henceforth we may assume $\sigma > \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$.

We use the same notation as in Lemma 3.5, especially (3.12). Then $|\tau_2'| < |\tau_1|$ by (3.15), and Lemma 2.5 yields $\mathtt{res} + \delta_1 = \tau_1 + \tau_2'$ [as in (3.12)], where the error $\delta_1$ is the second result of $\mathtt{FastTwoSum}(\tau_1, \tau_2')$. By (3.13) we know $s = \tau_1 + \tau_2' + \delta$, so that $\delta_1 \leq 0$ and (3.17) imply

$$s = \mathtt{res} + \delta_1 + \delta < \mathtt{res} + \frac{1}{2}\mathtt{eps}|\mathtt{res}| \leq \mathcal{M}^+(\mathtt{res}) .$$

Thus (7.2) shows $\mathrm{fl}(s) \neq \mathrm{succ}(\mathtt{res})$. The case $\delta_1 \geq 0$ follows similarly, and the case $\delta_1 = 0$ follows. $\qquad\square$

By Lemma 7.2 it can be decided which $\mu \in \{\mathcal{M}^-(\mathtt{res}), \mathcal{M}^+(\mathtt{res})\}$ has to be subtracted from $s$ to decide whether $\mathrm{fl}(s)$ is equal to $\mathtt{res}$ or to one of its neighbors. However, $\mu$ is not a floating-point number but could only be subtracted in two parts by splitting $\mu$ into $\mathtt{res}$ and $\delta_{\mathtt{res}} := \mu - \mathtt{res} \in \mathbb{F}$. Rather than extracting the vector $p$ appended by $-\mathtt{res}$ and $-\delta_{\mathtt{res}}$ again, it is much better to use the already extracted vector $p'$. Fortunately, our analysis in Lemma 6.3 gives an equation for $s - \mathtt{res}$ in terms of $p_i'$, namely $s - \mu = s - \mathtt{res} - \delta_{\mathtt{res}} = R + \sum p_i' - \delta_{\mathtt{res}}$ for $R$ as computed in Algorithm 6.2 ($\mathtt{TransformK}$). Now Algorithm 3.3 ($\mathtt{Transform}$) allows for one extra parameter to be added to $\sum p_i'$, but not two. The following lemma shows that one parameter $R' \in \mathbb{F}$ is sufficient and $\mathtt{Transform}$ can be applied directly to $p'$ and $R'$.

LEMMA 7.3. *Let $p$ be a vector of $n$ floating-point numbers, define $M := \lceil \log_2(n+2) \rceil$, and assume $2^{2M}\mathtt{eps} \leq 1$. Let $\mathtt{res}$, $R$ and $p'$ be the results of Algorithm 6.2 ($\mathtt{TransformK}$) applied to $p$ and $\varrho = 0$. Let $\delta_{\mathtt{res}} \in \{\frac{1}{2}(\mathrm{pred}(\mathtt{res}) - \mathtt{res}), \frac{1}{2}(\mathrm{succ}(\mathtt{res}) - \mathtt{res})\}$ and assume $\delta_{\mathtt{res}} \in \mathbb{F}$.*

*Then $R - \delta_{\mathtt{res}} \in \mathbb{F}$. Moreover, if the vector $p'$ is nonzero, then $\delta_{\mathtt{res}} \in \mathtt{eps}\sigma'\mathbb{Z}$ is satisfied for $\mu' := \max_i |p_i'|$ and $\sigma' := 2^{M + \lceil \log_2 \mu' \rceil}$.*

PROOF. For the computation of $\mathtt{res}$ we use the same code as in (3.7), the assumptions of Lemma 3.5 and Lemma 6.3 are satisfied and we may use results of the lemmas and the proofs, especially $\Delta = \mathtt{res} - \tau_1$ and $R = \tau_2 - \Delta$.

Assume the final value of $\sigma$ in $\mathtt{Transform}$ satisfies $\sigma \leq \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta}$. Then by Lemma 3.5 the vector $p'$ is entirely zero, $s = \tau_1 + \tau_2$ and $\mathtt{res} = \mathrm{fl}(\tau_1 + \tau_2) = \tau_1$. Thus $\Delta = \mathtt{res} - \tau_1 = 0$ and $R = \tau_2$. If the final value of $m$ in Algorithm 3.3 ($\mathtt{Transform}$) is $m = 1$, then $\max |p_i| \leq \mu \leq 2^{-M}\sigma$ so that $|s| \leq n\mu < \sigma$ implies $s \in \mathbb{U}$ and $\mathtt{res} = s \in \mathbb{U}$, which means $\delta_{\mathtt{res}} = \mathtt{eta}/2 \notin \mathbb{F}$, a contradiction. Therefore, the final value $m$ satisfies $m > 1$, and the "until"-condition in yields $|\mathtt{res}| = |\tau_1| = |t^{(m-1)}| < 2^{2M}\mathtt{eps}\sigma_{m-2} = 2^M\sigma \leq 2^{M-1}\mathtt{eps}^{-1}\mathtt{eta}$. Hence $|\delta_{\mathtt{res}}| \leq 2^{M-2}\mathtt{eta}$ and $|R| = |\tau_2| \leq \mathtt{eps}|\tau_1| \leq 2^{M-1}\mathtt{eta}$, and $R - \delta_{\mathtt{res}} \in \mathbb{F}$ follows.

Henceforth we assume $\sigma > \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta}$. Then the assumption $\delta_{\mathtt{res}} \in \mathbb{F}$ implies $|\mathtt{res}| \geq \mathtt{eps}^{-1}\mathtt{eta}$ and

$$(7.5) \qquad \delta_{\mathtt{res}} \quad \in \quad \begin{cases} \frac{1}{2}\mathtt{eps} \cdot \mathrm{ufp}(\mathtt{res})\mathbb{Z} & \text{if } |\mathtt{res}| = \mathrm{ufp}(\mathtt{res}) \\ \mathtt{eps} \cdot \mathrm{ufp}(\mathtt{res})\mathbb{Z} & \text{if } |\mathtt{res}| \neq \mathrm{ufp}(\mathtt{res}) \end{cases}$$

and

$$(7.6) \qquad |\delta_{\mathtt{res}}| \leq \mathtt{eps} \cdot \mathrm{ufp}(\mathtt{res}) .$$

We know by (3.16) and (3.11) that

$$|\mathtt{res}| > \frac{5}{8}|\tau_1| > 2^{2M-1}\mathtt{eps}\sigma \geq 2^{M+1}\mathtt{eps}\sigma .$$

Moreover (3.8) gives $\sigma' \leq 2^M\mathtt{eps}\sigma$, so that (7.5) and (2.4) prove $\delta_{\mathtt{res}} \in 2^M\mathtt{eps}^2\sigma\mathbb{Z} \subseteq \mathtt{eps}\sigma'\mathbb{Z}$.

To show $R - \delta_{\mathtt{res}} \in \mathbb{F}$, we distinguish two cases. First, assume $|\tau_1| < \sigma$. Then $|\tau_2| \leq \mathtt{eps}|\tau_1| < \mathtt{eps}\sigma$ and $\tau_2 \in \mathtt{eps}\sigma\mathbb{Z}$ as in (3.9) yields $\tau_2 = 0$. Hence $R \in \mathbb{F}$ gives $R = \tau_1 - \mathtt{res} = \mathrm{fl}(\tau_1 - \mathtt{res})$, and (2.8) shows $R \in \mathtt{eps}\cdot\mathrm{ufp}(\mathtt{res})\mathbb{Z} \subseteq \frac{1}{2}\mathtt{eps}\cdot\mathrm{ufp}(\mathtt{res})\mathbb{Z}$. Regarding (7.5) and (2.10) it suffices to show $|R - \delta_{\mathtt{res}}| \leq \frac{1}{2}\mathrm{ufp}(\mathtt{res})$. By (3.12), (2.9), $\tau_2 = 0$, (3.14) and (3.11),

$$(7.7) \quad \begin{aligned} |R| &= |\mathtt{res} - \tau_1| = |\tau_3 - \delta_2 - \delta_1| \leq |\tau_3| + \mathtt{eps}|\tau_3| + |\delta_1| \leq (1 + \mathtt{eps})n\mathtt{eps}\sigma + \mathtt{eps} \cdot \mathrm{ufp}(\mathtt{res}) \\ &\leq (2^M - 1)\mathtt{eps}\sigma + \mathtt{eps} \cdot \mathrm{ufp}(\mathtt{res}) \leq (2^M - 1)2^{-2M}\mathrm{ufp}(\tau_1) + \mathtt{eps} \cdot \mathrm{ufp}(\mathtt{res}) . \end{aligned}$$

Now (3.16) implies $\mathrm{ufp}(\tau_1) \le 2\mathrm{ufp}(\mathtt{res})$, so that (7.6) and the assumption $2^{2M}\mathtt{eps} \le 1$ yield

$$|R - \delta_{\mathtt{res}}| \le \big((2^M - 1)2^{-2M+1} + 2\mathtt{eps}\big)\mathrm{ufp}(\mathtt{res}) \le 2^{-M+1}\mathrm{ufp}(\mathtt{res}) \le \frac{1}{2}\mathrm{ufp}(\mathtt{res}) \ .$$

This finishes the first case. Second, assume $|\tau_1| \ge \sigma$. Then as in (10.7) we see

$$(7.8) \qquad\qquad |R| = |\tau_2 - \Delta| \le \mathtt{eps}(1 + 3\mathtt{eps})|\tau_1| + 2^M\mathtt{eps}\sigma \ ,$$

and by (2.8) we know $\Delta = \mathtt{res} - \tau_1 \in \mathtt{eps} \cdot \mathrm{ufp}(\tau_1) \subseteq \mathtt{eps}\sigma\mathbb{Z}$. We distinguish two subcases. First, suppose $|\mathtt{res}| \le \sigma$. Then (3.9) yields $R = \tau_2 - \Delta \in \mathtt{eps}\sigma\mathbb{Z} \subseteq \frac{1}{2}\mathtt{eps}\sigma\mathbb{Z}$. Regarding (7.5) implies $\delta_{\mathtt{res}} \in \frac{1}{2}\mathtt{eps} \cdot \mathrm{ufp}(\mathtt{res})\mathbb{Z} \subseteq \frac{1}{2}\mathtt{eps}\sigma\mathbb{Z}$, and with (2.10) it suffices to show $|R - \delta_{\mathtt{res}}| \le \frac{1}{2}\sigma$. This is seen with (7.8), (3.16), (7.6) and

$$|R - \delta_{\mathtt{res}}| \le \mathtt{eps}(1 + 3\mathtt{eps})\frac{8}{5}\sigma + 2^M\mathtt{eps}\sigma + \mathtt{eps}\sigma \le 3\mathtt{eps}\sigma + 2^{-M}\sigma < \frac{1}{2}\sigma.$$

This leaves us, secondly, with $|\tau_1| \ge \sigma$ and $|\mathtt{res}| > \sigma$. In this case (7.5) and (3.9) imply $\Delta, \delta_{\mathtt{res}} \in \mathtt{eps}\sigma\mathbb{Z}$ and $R = \tau_2 - \Delta \in \mathtt{eps}\sigma\mathbb{Z}$ as well, so with (2.10) it suffices to show $|R - \delta_{\mathtt{res}}| \le \sigma$. We have $t^{(0)} = \varrho = 0$ in Algorithm 3.3, so for the final value of $m$ in the "until"-condition ($\mathtt{Transform}$) and with (3.3) we have

$$|\tau_1| = |\mathrm{fl}(t^{(m-1)} + \tau^{(m)})| < (1 + \mathtt{eps})(2^{2M}\mathtt{eps}\sigma_{m-2} + \sigma) = (1 + \mathtt{eps})(2^M + 1)\sigma < (2^M + 2)\sigma \ .$$

Now (3.12) and (3.15) give

$$|\mathtt{res}| \le (1 + \mathtt{eps})|\tau_1 + \tau_2'| \le (1 + 2\mathtt{eps})^2|\tau_1| + 2^M\mathtt{eps}\sigma < (2^M + 3)\sigma \ ,$$

so that (7.6) yields $|\delta_{\mathtt{res}}| \le 2^M\mathtt{eps}\sigma$. Using (7.8) and $2^{2M}\mathtt{eps} \le 1$ we conclude

$$|R - \delta_{\mathtt{res}}| < (2^M + 4)\mathtt{eps}\sigma + 2^{M+1}\mathtt{eps}\sigma \le 2^{M+2}\mathtt{eps}\sigma \le \sigma \ .$$

The lemma is proved. □

We note that Lemma 7.3 does not remain true for $\varrho \ne 0$ as is seen by $\varrho = 1$, the 1-element vector $p = \mathtt{eps}^2$ and the choice $\delta_{\mathtt{res}} = -\mathtt{eps}$. However, we do not need this case in the sequel.

Now we can state the algorithm to compute the rounded-to-nearest result of a sum of floating-point numbers.

ALGORITHM 7.4. *Accurate summation with rounding to nearest.*

function $\mathtt{resN} = \mathtt{NearSum}(p)$
    $[\tau_1, \tau_2, p'] = \mathtt{Transform}(p, 0)$               % *Parameter* $\Phi$ *replaced by* $2^{2M}\mathtt{eps}$
    $\tau_2' = \mathrm{fl}(\tau_2 + (\sum_{i=1}^n p_i'))$
    $[\mathtt{res}, \delta] = \mathtt{FastTwoSum}(\tau_1, \tau_2')$           % $\mathtt{res} + \delta = \tau_1 + \tau_2'$
    if $\delta = 0$, $\mathtt{resN} = \mathtt{res}$, return, end if     % $\mathrm{fl}(s) = \mathtt{res}$
    $R = \tau_2 - (\mathtt{res} - \tau_1)$                 % $s - \mathtt{res} = R + \sum p_i'$
    if $\delta < 0$                        % $\mathrm{fl}(s) \in \{\mathrm{pred}(\mathtt{res}), \mathtt{res}\}$
        $\gamma = \mathrm{pred}(\mathtt{res}) - \mathtt{res}$           % $\mathtt{res} + \gamma = \mathrm{pred}(\mathtt{res})$
        if $\gamma = -\mathtt{eta}$, $\mathtt{resN} = \mathtt{res}$, return, end if   % $\mathcal{M}^-(\mathtt{res}) \notin \mathbb{F} \Rightarrow s = \mathtt{res}$
        $\delta' = \gamma/2$                     % $\mu := \mathtt{res} + \delta' = \mathcal{M}^-(\mathtt{res})$
        $\delta'' = \mathtt{TransformK}(p', R - \delta')$     % $s - \mu = R - \delta' + \sum p_i'$, $\delta'' \in \square(s - \mu)$
        if $\delta'' > 0$, $\mathtt{resN} = \mathtt{res}$         % $s > \mathcal{M}^-(\mathtt{res})$
            elseif $\delta'' < 0$, $\mathtt{resN} = \mathrm{pred}(\mathtt{res})$    % $s < \mathcal{M}^-(\mathtt{res})$
            else $\mathtt{resN} = \mathrm{fl}(\mathtt{res} + \delta')$        % $s = \mathcal{M}^-(\mathtt{res})$
        end if
    else                         % $\mathrm{fl}(s) \in \{\mathtt{res}, \mathrm{succ}(\mathtt{res})\}$
        the case $\delta > 0$ is treated similarly
    end if

THEOREM 7.5. *Let $p$ be a vector of $n$ floating-point numbers, define $M := \lceil \log_2(n+2) \rceil$, and assume $2^{2M}\mathsf{eps} \le 1$. Let $\mathsf{resN}$ be the result of Algorithm 7.4 (NearSum) applied to $p$.*

*Then $\mathsf{resN} = \mathrm{fl}(s)$ is the rounded-to-nearest exact sum $s := \sum_{i=1}^n p_i$ in the sense of IEEE 754.*

REMARK 1. When the floating-point rounding "fl" is omitted in Algorithm 7.4, then we will show that no rounding error can occur.

REMARK 2. Note that numerical evidence suggests that it seems worth to check for $\delta = 0$ to save the second extraction by TransformK. If the "repeat-until"-loop in the first and second extraction by TransformK is executed $m$ and $m'$ times, respectively, and using the fact that the second extraction unit $\sigma$ needs not to be computed, then NearSum needs $(4m + 4m' + 4)n + \mathcal{O}(m + m')$ flops.

REMARK 3. Note that we generally assume that no overflow occurs, however, this is easily treated by some scaling.

REMARK 4. We used the predecessor and successor of a floating-point number. These are especially easy to calculate if directed rounding as in IEEE 754 is available. Otherwise, for example, they may computed in rounding to nearest using Algorithm 3.5 (NextPowerTwo) from Part I of this paper, or by the algorithms presented in [23].

PROOF OF THEOREM 7.5. The internal result $\mathsf{res}$ is computed in the same way as in (7.3), and the assumptions of Lemma 7.2 are satisfied. The quantity $\delta_1$ in Lemma 7.2 is the same as $\delta$ in NearSum. Therefore, $\mathrm{fl}(s) \in \{\mathrm{pred}(\mathsf{res}), \mathsf{res}, \mathrm{succ}(\mathsf{res})\}$, and $\delta = 0$ implies $\mathrm{fl}(s) = \mathsf{res}$. The quantity $R$ is computed exactly the same way as in Algorithm 6.2 (TransformK), so (6.1) implies $s - \mathsf{res} = R + \sum p_i'$.

If $\delta < 0$, then (7.4) implies $\mathrm{fl}(s) \in \{\mathrm{pred}(\mathsf{res}), \mathsf{res}\}$. The quantity $\gamma = \mathrm{pred}(\mathsf{res}) - \mathsf{res}$ is always in $\mathbb{F}$, and $\gamma = -\mathsf{eta}$ together with $\mathsf{res} \in \square(s)$ and $s = \sum p_i \in \mathsf{eta}\mathbb{Z}$ implies $s = \mathsf{res}$. If $\gamma \ne -\mathsf{eta}$, then $\delta' = \frac{1}{2}(\mathrm{pred}(\mathsf{res}) - \mathsf{res}) \in \mathbb{F}$ and $\mathcal{M}^-(\mathsf{res}) = \mathsf{res} + \delta'$. The quantity $\delta_{\mathsf{res}}$ in Lemma 7.3 is the same as $\delta'$, so $R - \delta' \in \mathbb{F}$. Moreover, using $\sigma' = 2^{M + \lceil \log_2 \mu' \rceil}$ with $\mu' := \max_i |p_i'|$ as in Lemma 6.3, (6.2) and Lemma 7.3 prove $R, \delta', R - \delta' \in \mathsf{eps}\sigma'\mathbb{Z}$ for the first extraction of $(p, 0)$, so that the assumptions of Lemma 6.3 for the application of Algorithm 6.2 (TransformK) to $(p', R - \delta')$ are satisfied. Therefore $\delta'' \in \square(R - \delta' + \sum p_i')$, where $R - \delta' + \sum p_i' = s - \mathsf{res} - \delta' = s - \mathcal{M}^-(\mathsf{res})$. Hence $\mathrm{sign}(s - \mathcal{M}^-(\mathsf{res})) = \mathrm{sign}(\delta'')$. So the assertion follows by (7.2) if $\delta'' \ne 0$. Finally, if $\delta'' = 0$, then $s = \mathcal{M}^-(\mathsf{res})$ and $\mathrm{fl}(s) = \mathrm{fl}(\mathcal{M}^-(\mathsf{res}))$.

The case $\delta > 0$ is programmed and treated similarly, and the theorem is proved. $\square$

In contrast to Algorithm 4.5 (AccSum) in Part I, the computing time of Algorithm 7.4 (NearSum) depends on the exponent variation of the vector entries rather than on the condition number of the sum. However, it seems unlikely that the maximum number of 40 extractions for IEEE 754 double precision is achieved in other than constructed examples.

There is an apparent contradiction, namely that the computing time of AccSum is proportional to the logarithm of the condition number, but that of NearSum with only one extra call of Algorithm 6.2 (TransformK) is not. However, the computing time of AccSum to compute a faithfully rounded result is proportional to $t_f := \log(\mathrm{cond}(\sum p_i))$, but that of NearSum is proportional to $t_N := \log(\mathrm{cond}(R - \delta' + \sum p_i'))$, where $R - \delta' + \sum p_i'$ is the difference between the exact sum $s$ and one of the "switching points" $\mathcal{M}^-(\mathsf{res})$ or $\mathcal{M}^+(\mathsf{res})$. Note that, subject to the size of the exponent range, $t_N / t_f$ can be arbitrarily large.

Finally we mention that combining the results of Sections 6 and 7 we can easily define an algorithm producing a sequence of $K$ floating-point numbers representing a $K$-fold rounded-to-nearest result of $s = \sum p_i$. Consider Algorithm 6.4 (AccSumK) changed in such a way that only the last member $\mathsf{Res}_K$ of the sequence $(\mathsf{Res}_1, \cdots, \mathsf{Res}_k)$ is computed in rounding-to-nearest using a piece of code similar to that in Algorithm 7.4

(NearSum). Then (5.14) implies in the notation of the proof of Proposition 5.6

$$\left| s - \sum_{k=1}^{K} \texttt{Res}_k \right| = |\Delta_K| \leq \texttt{eps} \cdot \text{ufp}(\Delta_{K-1}) \leq \texttt{eps}^K \text{ufp}(s) \leq \texttt{eps}^K |s| \ .$$

Hence $\sum \texttt{Res}_k$ is of the same accuracy as a nearest floating-point approximation in a floating-point grid with relative rounding error unit $\texttt{eps}^K$, and the last member $\texttt{Res}_K$ is $s - \sum_{k=1}^{K-1} \texttt{Res}_k$ rounded-to-nearest.

**8. 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 restriction imposed by the assumption $2^{2M}\texttt{eps} \leq 1$. This should suffice for most practical purposes. In IEEE 754 single precision with $\texttt{eps} = 2^{-24}$, the number of summands is restricted to $n = 4094$, which may be an obstacle for the application of our algorithms.

In Section 4 we showed that using $2^M\texttt{eps}$ for the parameter $\Phi$ in Algorithm 3.3 (Transform) we have enough information to calculate the sign of a sum. This was also true for huge vector lengths. Next we show how to compute a faithfully rounded result of the sum for huge vector lengths. For this we continue to extract $p_i'$ by Transform until the error term of the summation of the extracted vector is small enough compared to $\tau_1$. This is done in the following Algorithm 8.1 (AccSumHugeN).

ALGORITHM 8.1. *Accurate summation with faithful rounding for huge $n$.*

> function $\texttt{res} = \texttt{AccSumHugeN}(p)$
>     $[\tau_1, \tau_2, q^{(0)}, \sigma, M] = \texttt{Transform}(p, 0)$             % *Parameter $\Phi$ replaced by $2^{M+3}\texttt{eps}$*
>     if $\sigma \leq \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$, $\texttt{res} = \tau_1$, return, end if
>     $k = 0$; $\sigma_0 = \text{fl}(2^M\texttt{eps}\sigma)$
>     repeat
>         $k = k + 1$
>         $[\tau^{(k)}, q^{(k)}] = \texttt{ExtractVector}(\sigma_{k-1}, q^{(k-1)})$
>         $\sigma_k = \text{fl}(2^M\texttt{eps}\sigma_{k-1})$
>     until  $|\tau_1| \geq \text{fl}\big((2^{2M+1}\texttt{eps})\sigma_{k-1}\big)$    or    $\sigma_{k-1} \leq \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$
>     $\tilde{\tau}^{(1)} = \text{fl}\big(\tau^{(1)} + (\tau^{(2)} + \cdots + (\tau^{(k)} + (\tau_2 + (\sum_{i=1}^{n} q_i^{(k)}))) \cdots)\big)$
>     $\texttt{res} = \text{fl}(\tau_1 + \tilde{\tau}^{(1)})$

PROPOSITION 8.2. *Let $p$ be a vector of $n$ floating-point numbers, assume $2^{M+3}\texttt{eps} \leq 1 < 2^{2M}\texttt{eps}$ and $\texttt{eps} \leq 1/512$. Let $\texttt{res}$ be the result of Algorithm 8.1 (AccSumHugeN) applied to $p$.*

*Then $\texttt{res}$ is a faithful rounding of $s := \sum_{i=1}^{n} p_i$. The algorithm needs $(4m + 4K + 3)n + \mathcal{O}(m + K)$ flops for $m$ "repeat-until"-loops within Transform and $K$ denoting the final value of $k$.*

REMARK 1. For IEEE 754 single precision with $\texttt{eps} = 2^{-24}$ this limits the vector length to a little over 2 million rather than $n = 4094$ for Algorithm AccSum. For IEEE 754 double precision the vector length is now limited to about $1.1 \cdot 10^{15}$.

REMARK 2. For huge vector lengths near the admissible maximum, however, the algorithm becomes rather inefficient because few bits, possibly only three bits, are extracted in each execution of Transform.

REMARK 3. We mention that the parameters can be adjusted so that the weaker assumption $2^{M+2}\texttt{eps} \leq 1$ suffices. However, this extracts ultimately only two bits at a time.

REMARK 4. The assumption $1 < 2^{2M}\texttt{eps}$ imposes no restriction because otherwise Algorithm 4.5 (AccSum) in Part I can (and for better performance should) be used.

PROOF OF PROPOSITION 8.2. Algorithm `Transform` is called with $\varrho = 0$. So $2^M \mathtt{eps} \leq \frac{1}{8}$ and Lemma 3.4 imply

$$(8.1) \qquad\qquad s = \tau_1 + \tau_2 + \sum_{i=1}^{n} q_i^{(0)} \ ,$$

$$(8.2) \qquad\qquad \max |q_i^{(0)}| \leq \mathtt{eps}\sigma,$$

$$(8.3) \qquad\qquad \tau_1 = \mathrm{fl}(\tau_1 + \tau_2) \quad \text{and} \quad |\tau_2| \leq \mathtt{eps}|\tau_1| \ .$$

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

$$(8.4) \qquad \varphi = 2^M \mathtt{eps} \leq \frac{1}{8} \ , \quad \mathtt{eps} \leq \frac{1}{512} \ , \quad M \geq 5 \quad \text{and} \quad \sigma_k = \varphi^{k+1}\sigma \ \text{for} \ 0 \leq k \leq K \ ,$$

and the choice of $\Phi$ and (3.5) gives

$$(8.5) \qquad\qquad |\tau_1| \geq 8\varphi\sigma \ .$$

We first show that we may assume

$$(8.6) \qquad\qquad \sigma > \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta} \qquad \text{and} \qquad |\tau_1| \geq \mathtt{eps}^{-1}\mathtt{eta}$$

without loss of generality. To see that, first assume that $\sigma$, computed by Algorithm `Transform`, is in the underflow range. Then $q^{(0)}$ is the zero vector by Lemma 3.4, and $\mathrm{fl}(s) = \mathtt{res}$ by (8.1). Second, assume $|\tau_1| < \mathtt{eps}^{-1}\mathtt{eta}$. Then (8.5) implies $8\varphi\sigma \leq |\tau_1| < \mathtt{eps}^{-1}\mathtt{eta}$, so that by $\sigma_0 = \varphi\sigma < \frac{1}{8}\mathtt{eps}^{-1}\mathtt{eta}$ the "repeat-until"-loop is finished for $k = 1$. Moreover $|\tau_2| \leq \mathtt{eps} \cdot \mathrm{ufp}(\tau_1) \leq \frac{1}{2}\mathtt{eta}$ implies $\tau_2 = 0$, and by Theorem 3.2, $|q_i^{(1)}| \leq \mathtt{eps}\sigma_0$, so that the vector $q_i^{(1)}$ is identically zero. Therefore $\mathtt{res} = \mathrm{fl}(\tau_1 + \tau^{(1)}) = \mathrm{fl}(s)$ by (8.7), and $\mathtt{res}$ is again a faithful rounding of $s$.

By definition of $\sigma_0$ and (8.2), $\max_i |q_i^{(0)}| \leq 2^{-M}\sigma_0$, so the assumptions of Theorem 3.2 are satisfied for the first call of `ExtractVector` with $k = 1$. Therefore,

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

If $\sigma_0 \leq \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta}$, 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 $\sigma_0 > \frac{1}{2}\mathtt{eps}^{-1}\mathtt{eta}$, then $\sigma_1$ is computed without rounding error and can safely be used. Note that the assumption $1 < 2^{2M}\mathtt{eps}$ implies that no rounding error occurs in the computation of $\mathrm{fl}\big((2^{2M+1}\mathtt{eps})\sigma_{k-1}\big)$ in the "until"-condition.

Again, the assumptions of Theorem 3.2 are satisfied for the next call of `ExtractVector`, and repeating this argument shows

$$(8.7) \qquad s = \tau_1 + \tau_2 + \sum_{\nu=1}^{k} \tau^{(\nu)} + \sum_{i=1}^{n} q_i^{(k)} \ , \quad |\tau^{(k)}| < \sigma_{k-1} \quad \text{and} \quad \max_i |q_i^{(k)}| \leq 2^{-M}\sigma_k$$

for $k$ between 1 and its final value $K$. The same argument as before applies to possible rounding errors in the computation of $\sigma_K$. Denote

$$(8.8) \qquad \tau^{(K+1)} := \tau_2 + \big(\sum_{i=1}^{n} q_i^{(K)}\big) \quad \text{and} \quad \tilde{\tau}^{(K+1)} := \mathrm{fl}\big(\tau_2 + (\sum_{i=1}^{n} q_i^{(K)})\big) = \tau^{(K+1)} - \delta_{K+1},$$

$$(8.9) \qquad \tilde{\tau}^{(k)} := \mathrm{fl}\big(\tau^{(k)} + \tilde{\tau}^{(k+1)}\big) = \tau^{(k)} + \tilde{\tau}^{(k+1)} - \delta_k \quad \text{for} \ 1 \leq k \leq K \ .$$

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

$$(8.10) \qquad \tilde{\tau}^{(1)} = \sum_{k=1}^{K+1} \tau^{(k)} - \sum_{k=1}^{K+1} \delta_k = \sum_{k=1}^{K} \tau^{(k)} + \tau_2 + \sum_{i=1}^{n} q_i^{(K)} - \sum_{k=1}^{K+1} \delta_k \ .$$

Denote $\tau' := \sum q_i^{(K)}$ and $\tilde{\tau}' := \mathrm{fl}\big(\sum q_i^{(K)}\big)$. Then by definition (8.8), (8.7), (2.7) and $n + 2 \leq 2^M$,

$$
\begin{aligned}
(8.11) \qquad |\delta_{K+1}| \;&=\; |\tilde{\tau}^{(K+1)} - \tau^{(K+1)}| = |\mathrm{fl}(\tau_2 + \tilde{\tau}') - (\tau_2 + \tilde{\tau}') + (\tau_2 + \tilde{\tau}') - (\tau_2 + \tau')| \\
&\leq\; \mathsf{eps}|\tau_2 + \mathrm{fl}\big(\textstyle\sum q_i^{(K)}\big)| + |\mathrm{fl}\big(\textstyle\sum q_i^{(K)}\big) - \textstyle\sum q_i^{(K)}| \\
&\leq\; \mathsf{eps}\big(|\tau_2| + n2^{-M}\sigma_K + \tfrac{1}{2}n(n-1)2^{-M}\sigma_K\big) \\
&<\; \mathsf{eps}\big(|\tau_2| + 2^{M-1}\sigma_K\big) \;.
\end{aligned}
$$

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

$$
(8.12) \qquad |\tilde{\tau}^{(k)}| \leq \Big( \sum_{\nu=0}^{K+1-k} (1+\mathsf{eps})^{\nu+1}\varphi^{\nu}\Big)\sigma_{k-1} + (1+\mathsf{eps})^{K+2-k}|\tau_2| \quad \text{for } 1 \leq k \leq K+1 \;.
$$

For $k = K+1$ this follows by (8.5)

$$
|\tilde{\tau}^{(K+1)}| \leq (1+\mathsf{eps})|\tau_2 + \mathrm{fl}\big(\textstyle\sum q_i^{(K)}\big)| \leq (1+\mathsf{eps})(|\tau_2| + n2^{-M}\sigma_K) < (1+\mathsf{eps})(|\tau_2| + \sigma_K) \;,
$$

again using (2.7), and the induction step uses

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

by (8.9), and $|\tau^{(k)}| < \sigma_{k-1}$ by (8.7). Hence (8.4) gives

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

and by (8.9) and (8.4) it follows

$$
(8.13) \qquad \mathsf{eps}^{-1}|\delta_k| \leq |\tilde{\tau}^{(k)}| \leq \frac{147}{128}\varphi^k\sigma + (1+\mathsf{eps})^{K+2-k}|\tau_2| \quad \text{for } 1 \leq k \leq K \;.
$$

Next we need an upper bound on $K$, the final value of $k$. Denote $\mathsf{eps} = 2^{-m}$, then (8.4) gives $m \geq 9$ and $M - m \leq -3$. If $K \geq 2$, then the "until"-condition, (8.5) and (8.4) imply $8\varphi\sigma \leq |\tau_1| < 2^{2M+1}\mathsf{eps}\sigma_{K-2} = 2^{M+1}\sigma_{K-1} = 2^{M+1}\varphi^K\sigma$, so that $8\varphi \leq 2^M\varphi^K$. Hence $K - 1 \leq \frac{M-3}{m-M}$, so that

$$
(8.14) \qquad\qquad K \leq \frac{M}{3} \quad \text{and} \quad K \leq \frac{m-3}{3} \;,
$$

which is by (8.4) also true for $K < 2$. Therefore, $m \geq 9$ yields

$$
(8.15) \qquad (1+\mathsf{eps})^K \leq (1+2^{-m})^{\frac{m-3}{3}} < \frac{129}{128} \quad \text{and} \quad (1+\mathsf{eps})^2\frac{(1+\mathsf{eps})^K - 1}{\mathsf{eps}} \leq \frac{1}{128}\mathsf{eps}^{-1} \;.
$$

Now the definition of $\mathtt{res}$, (8.13), (8.3), (8.15), (8.4) and (8.5) yield

$$
\begin{aligned}
(8.16) \qquad |\mathtt{res}| \;&\geq\; (1-\mathsf{eps})\big(|\tau_1| - |\tilde{\tau}^{(1)}|\big) \\
&\geq\; (1-\mathsf{eps})\big(|\tau_1| - \tfrac{147}{128}\varphi\sigma - (1+\mathsf{eps})^{K+1}\mathsf{eps}|\tau_1|\big) \\
&\geq\; \tfrac{127}{128}|\tau_1| - \tfrac{147}{128}\varphi\sigma \;.
\end{aligned}
$$

Hence (8.5) yields

$$
(8.17) \qquad\qquad |\mathtt{res}| > \frac{217}{256}|\tau_1| \;,
$$

and (8.6) implies $\mathtt{res} \notin \mathbb{U}$. Now set $r := \tau_1 + \tilde{\tau}^{(1)}$, so that $\mathtt{res} = \mathrm{fl}(r)$. Then (8.7) and (8.10) imply

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

| $K$ | single precision | double precision |
|:---:|:---:|:---:|
| 1 | 8190 | $1.3 \cdot 10^8$ |
| 2 | 65534 | $6.8 \cdot 10^{10}$ |
| 3 | $2.6 \cdot 10^5$ | $1.0 \cdot 10^{12}$ |
| 4 | $5.2 \cdot 10^5$ | $4.3 \cdot 10^{12}$ |
| 5 | $1.0 \cdot 10^6$ | $1.7 \cdot 10^{13}$ |

We will show $2|\delta| < \mathtt{eps}|\mathtt{res}|$ and apply Lemma 2.3 to demonstrate that $\mathtt{res}$ is a faithful rounding of the true result $s$. First we need to bound $\delta_{K+1}$. If $\sigma_{K-1} \in \mathbb{U}$, then (8.7) implies $q_i^{(K)} = 0$ for all $i$, and (8.8) gives $\delta_{K+1} = 0$. Otherwise the "until"-condition and (8.11) yield

$$2^{M-1}\sigma_K = 2^{2M-1}\mathtt{eps}\sigma_{K-1} \le \frac{1}{4}|\tau_1| \quad \text{and} \quad \mathtt{eps}^{-1}|\delta_{K+1}| \le (\mathtt{eps} + \frac{1}{4})|\tau_1| \ .$$

Hence (8.13), (8.15), $\mathtt{eps} \le \frac{1}{512}$ and $\varphi\sigma \le \frac{1}{8}|\tau_1|$ show

(8.18)
$$\begin{aligned}
\mathtt{eps}^{-1}|\delta| \quad &< \quad \frac{147}{128}\frac{\varphi}{1-\varphi}\sigma + (1+\mathtt{eps})^2 \frac{(1+\mathtt{eps})^K - 1}{\mathtt{eps}}|\tau_2| + \mathtt{eps}^{-1}|\delta_{K+1}| \\
&\le \quad \frac{147}{112}\varphi\sigma + \left(\frac{1}{128} + \mathtt{eps} + \frac{1}{4}\right)|\tau_1| \\
&\le \quad \frac{217}{512}|\tau_1| \ ,
\end{aligned}$$

and (8.17) gives $2|\delta| < \mathtt{eps}|\mathtt{res}|$. Therefore $\mathtt{res} \notin \mathbb{U}$ and Lemma 2.3 prove that $\mathtt{res}$ is a faithful rounding of $s$. The proof is finished. $\qquad\square$

With (8.14) we see that for IEEE 754 double precision at most $K = 17$ extractions are possible, and we can also estimate the minimum treatable vector length $n$ by Algorithm 8.1 (`AccSumHugeN`) depending on the number of extra extractions. Suppose the "repeat-until"-loop in Algorithm 8.1 (`AccSumHugeN`) has been executed $K$ times, and assume

(8.19)
$$M \le \frac{mK + 2}{K + 1} \ .$$

Then the first call of `Transform` implies $|\tau_1| \ge 8\varphi\sigma$ and a little computation using $\mathtt{eps} = 2^m$ yields

$$|\tau_1| \ge 2^{M+3-m} \ge 2^{2M+1}\mathtt{eps}\varphi^K \ ,$$

so that the "until"-condition is satisfied. That means, if $M$ satisfies (8.19), then the "repeat-until"-loop in `AccSumHugeN` is executed at most $K$ times. It follows the *minimum* treatable vector lengths $n$ with $K$ loops, which are summarized in Table 8.1. As expected, `AccSumHugeN` becomes inefficient for a vector length approaching $\mathtt{eps}^{-1}$; in IEEE 754 single precision about $10^6$ and in double precision about $2 \cdot 10^{13}$ may be a reasonable limit for $n$. As before an algorithm with $K$-fold faithful rounding and rounding-to-nearest for input vectors of huge length may be developed as well.

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

All algorithms are tested in three different environments, namely Intel Pentium 4, Intel Itanium 2 and AMD Athlon 64. We carefully choose compiler options to achieve best possible results, see Table 9.1.

We faced no problems except for Pentium 4 and the Intel Visual Fortran 9.1 compiler, where the code optimization/simplification is overdone by the compiler. A typical example is the first line $q_i = \mathrm{fl}\left((\sigma + p_i) - \sigma\right)$ in Algorithm 3.1 (`ExtractVector`), which is optimized into $q_i = p_i$. This can, of course, be avoided by setting appropriate compiler options; however, this may slow down the whole computation. In this specific case the

TABLE 9.1
*Testing environments*

| | CPU, Cache sizes | Compiler, Compile options |
| | Peak performance (for summation) | Memory bandwidth (approx.) |
|---|---|---|
| I) | Intel Pentium 4 (2.53GHz) | Intel Visual Fortran 9.1 |
| | L2: 512KB | /O3 /QaxN /QxN  [/Op, see Table 9.2] |
| | 2.53GFlops | 1GB/s |
| II) | Intel Itanium 2 (1.4GHz) | Intel Fortran 9.0 |
| | L2: 256KB, L3: 3MB | -O3 |
| | 2.8GFlops | 2.6GB/s |
| III) | AMD Athlon 64 (2.2GHz) | GNU gfortran 4.1.1 |
| | L2: 512KB | -O3 -fomit-frame-pointer |
| | | -march=athlon64 -funroll-loops |
| | 2.2GFlops | 3GB/s |

TABLE 9.2
*Compile options for Pentium 4, Intel Visual Fortran 9.1*

| Algorithm | Necessity of compile option /Op |
|---|---|
| SSum | No |
| Sum2 | Yes, for TwoSum |
| XBLAS | Yes, for TwoSum and FastTwoSum |
| Priest, PriestS | Yes, for FastTwoSum |
| Malcolm, MalcK, MalcN, MalcS, LAccu, LAccuK, LAccuN, LAccuS | Yes, for Split |
| AccSum, AccSumK, NearSum, AccSumHugeN, AccSign | Basically, no |

second author suggested a simple trick to overcome this by using $q_i = \mathrm{fl}\,(|\sigma + p_i| - \sigma)$ instead. This does not change the intended result since $|p_i| \leq \sigma$ is assumed in the analysis (Theorem 3.2), it avoids unintended compiler optimization, and it does not slow down the computation. For the other algorithms to be tested we had to use, however, the compile option /Op for Pentium 4. This ensures the consistency of IEEE standard 754 floating-point arithmetic. The compile options for the different algorithms are summarized in Table 9.2.

To test the algorithms presented in this paper, examples for huge condition numbers larger than $\mathtt{eps}^{-1}$ were generated by Algorithm 6.1 in [18], where a method to generate a vector whose summation is arbitrarily ill-conditioned is described. All tests are performed in IEEE 754 double precision except those for AccSumHugeN, which are performed in IEEE 754 single precision.

First, test results for Algorithm 6.4 (AccSumK) for $K$-fold faithfully rounded results are presented. For all examples we choose dimension $n = 1000$ and small condition number around $10^3$ or so and watch the effect of increasing values of $K$. If the condition number of the initial sum and of sums of intermediately extracted vectors does not exceed $\mathtt{eps}^{-1}$ significantly, we may expect the computing time to grow linearly in $K$.

Competitors for $K$-fold rounding are Malcolm's summation [16] and the long accumulator [12]. In fact, both approaches produce the bit representation of the exact result in an intermediate step: Malcolm's in an overlapping sequence of floating-point numbers, and the long accumulator in a non-overlapping sequence. From this it is not too difficult to extract a result of $K$-fold accuracy; the corresponding algorithms are denoted by MalcK and LAccuK, respectively. In Table 9.3 we normed the computing time of AccSum to 1, practically the same computing time as for AccSumK for $K = 1$.

Indeed we observe the linear dependency of AccSumK in $K$. For MalcK and LAccuK there is not too much dependency on $K$ since the exact result of the sum is computed anyway; only some additional effort to

TABLE 9.3

*Measured computing times for* AccSumK, $n = 1000$, *cond* $\sim 10^3$, *for all environments time of* AccSum *normed to 1*

| CPU | Intel Pentium 4 (2.53GHz) | | | Intel Itanium 2 (1.4GHz) | | | AMD Athlon 64 (2.2GHz) | | |
|---|---|---|---|---|---|---|---|---|---|
| Compiler | Intel Visual Fortran 9.1 | | | Intel Fortran 9.0 | | | GNU gfortran 4.1.1 | | |
| $K$ | MalcK | LAccuK | AccSumK | MalcK | LAccuK | AccSumK | MalcK | LAccuK | AccSumK |
| 1 | 17.5 | 91.7 | 1.0 | 14.4 | 54.1 | 1.0 | 6.0 | 28.0 | 1.0 |
| 2 | 17.1 | 89.8 | 1.6 | 15.5 | 57.4 | 1.4 | 6.1 | 28.8 | 1.8 |
| 3 | 18.4 | 89.8 | 2.2 | 15.5 | 57.1 | 1.8 | 6.2 | 28.4 | 3.1 |
| 4 | 19.2 | 89.8 | 3.4 | 15.9 | 57.4 | 2.5 | 6.3 | 28.6 | 3.9 |
| 5 | 20.0 | 91.8 | 3.9 | 16.2 | 57.6 | 2.9 | 6.5 | 29.1 | 4.8 |
| 6 | 22.4 | 89.8 | 4.6 | 16.5 | 57.5 | 3.4 | 6.7 | 28.9 | 5.6 |
| 7 | 23.7 | 91.8 | 5.7 | 16.9 | 57.7 | 4.1 | 6.9 | 28.8 | 6.9 |
| 8 | 27.5 | 95.8 | 6.4 | 17.4 | 57.9 | 4.5 | 7.2 | 28.9 | 7.7 |
| 9 | 30.0 | 93.8 | 7.0 | 17.9 | 58.1 | 4.9 | 7.6 | 29.3 | 8.5 |
| 10 | 32.2 | 93.9 | 8.0 | 18.5 | 58.4 | 5.6 | 7.9 | 29.5 | 9.8 |

TABLE 9.4

*Measured computing times for* NearSum, $n = 1000$, *for all environments time of* AccSum *normed to 1*

| CPU | Intel Pentium 4 (2.53GHz) | | | Intel Itanium 2 (1.4GHz) | | | AMD Athlon 64 (2.2GHz) | | |
|---|---|---|---|---|---|---|---|---|---|
| Compiler | Intel Visual Fortran 9.1 | | | Intel Fortran 9.0 | | | GNU gfortran 4.1.1 | | |
| $\text{cond}_{\text{near}}$ | MalcN | LAccuN | NearSum | MalcN | LAccuN | NearSum | MalcN | LAccuN | NearSum |
| $10^8$ | 17.3 | 92.2 | 1.7 | 14.7 | 57.2 | 1.7 | 5.5 | 28.3 | 1.8 |
| $10^{16}$ | 16.9 | 92.3 | 1.7 | 14.7 | 58.0 | 1.7 | 5.5 | 27.5 | 1.8 |
| $10^{24}$ | 16.5 | 92.3 | 2.2 | 14.9 | 58.4 | 2.0 | 5.7 | 28.2 | 2.3 |
| $10^{32}$ | 17.3 | 94.1 | 2.2 | 14.9 | 58.4 | 2.0 | 5.7 | 27.9 | 2.3 |
| $10^{40}$ | 16.6 | 90.6 | 2.7 | 15.4 | 58.8 | 2.5 | 5.9 | 29.0 | 2.8 |
| $10^{48}$ | 17.3 | 92.3 | 3.2 | 15.4 | 59.0 | 2.6 | 5.9 | 28.3 | 3.2 |
| $10^{56}$ | 17.3 | 94.2 | 3.1 | 15.5 | 58.9 | 2.7 | 6.1 | 29.1 | 3.2 |
| $10^{64}$ | 17.7 | 92.3 | 3.7 | 15.7 | 61.7 | 2.9 | 6.4 | 29.1 | 3.7 |

extract from this a sequence of $K$ floating-point numbers is necessary. On AMD Athlon 64 and for large $K$, Malcolm's approach is superior to AccSumK. Note that the accuracy of the results $\sum \text{Res}_k$ may in fact be better than $K$-fold since there may be sequences of adjacent zeros in the bit representation of the exact result, producing gaps between adjacent $\text{Res}_k$.

Next we tested NearSum. Challenging examples for rounding to nearest have an exact sum near the midpoint of two adjacent floating-point number. The usual condition number for the sum $\sum p_i$ is $\sum |p_i|/|\sum p_i|$, cf. (4.25) in Part I of this paper, measuring the amount of cancelation. For rounding to nearest we may use

$$(9.1) \qquad \text{cond}_{\text{near}} := \frac{\sum |p_i|}{\min\{|\frac{1}{2}(f + \text{succ}(f)) - \sum p_i| : f \in \mathbb{F}\}} \, ,$$

measuring the nearness of the sum to a "switching-point" for rounding to nearest. In Table 9.4 we normed again the computing time for AccSum to 1, so the additional amount of computing time to go from faithful rounding to rounding to nearest is monitored. Again, Malcolm's algorithm and the long accumulator compute a bit representation of the exact sum, and from there it is not difficult to derive the rounded-to-nearest result. The corresponding algorithms are denoted by MalcN and LAccuN, respectively. In all examples we choose the dimension $n = 1000$.

For small condition number the ratio of computing times for NearSum shows the expected factor 2 compared

Table 9.5

*Measured computing times for* `AccSumHugeN` *in single precision, cond* $= 10^8$*, time of* `SSum` *normed to 1*

| CPU | Intel Pentium 4 (2.53GHz) | | | | | Intel Itanium 2 (1.4GHz) | | | | |
| Compiler | Intel Visual Fortran 9.1 | | | | | Intel Fortran 9.0 | | | | |
| $n$ | Sum2 | XBLAS | Malcolm | LAccu | AccSum | Sum2 | XBLAS | Malcolm | LAccu | AccSum |
|---|---|---|---|---|---|---|---|---|---|---|
| 10,000 | 40.0 | 150.0 | 888.9 | 777.8 | 25.6 | 7.9 | 55.9 | 858.3 | 930.5 | 32.5 |
| 20,000 | 32.7 | 122.7 | 772.7 | 681.8 | 20.0 | 7.9 | 55.6 | 900.0 | 940.5 | 32.8 |
| 40,000 | 32.7 | 122.7 | 681.8 | 681.8 | 23.6 | 7.9 | 55.6 | 867.6 | 940.5 | 36.9 |
| 80,000 | 30.8 | 108.3 | 911.5 | 607.6 | 45.8 | 6.9 | 48.3 | 1058.8 | 816.7 | 35.8 |
| 160,000 | 4.5 | 15.8 | 125.5 | 94.1 | 21.6 | 6.4 | 44.4 | 961.5 | 743.4 | 36.1 |
| 320,000 | 3.6 | 12.7 | 176.9 | 70.8 | 23.3 | 6.0 | 40.4 | 1398.7 | 682.0 | 48.3 |
| 640,000 | 3.6 | 13.3 | 176.1 | 76.6 | 29.3 | 4.3 | 28.4 | 965.1 | 472.0 | 53.4 |
| 1,280,000 | 3.5 | 13.0 | 242.3 | 71.9 | 36.9 | 4.2 | 26.5 | 1253.2 | 445.1 | 63.0 |

to `AccSum` caused by the additional call of `TransformK`. Also as expected, the computing time of `MalcN` and `LAccuN` is almost independent of the condition number. For $n = 1000$ some 43 bits are extracted in double precision at a time. Hence for condition number $10^{64}$ which is about $2^{213}$ we need some 5 extractions, and this is reflected in the computing time ratio of `NearSum` to `AccSum` in Table 9.4. Note that condition numbers exceeding $10^{16}$ occur only in very special applications.

In Figures 5.1 and 5.3 in Part I of this paper we displayed the MFlop-rates for different algorithms. They showed that `AccSum` achieves a much better MFlop-rate than Malcolm's or long accumulator algorithm. For `AccSumK` and `NearSum` this is similar; for small condition numbers they achieve about 85% to 95%, for larger ones between 105% up to 160% and a little more of the MFlop-rate of `AccSum`. This corresponds to about 50% to 80% of the peak performance.

Next we tested `AccSumHugeN`. For double precision we may use `AccSum` until dimension $n = 6.4 \cdot 10^7$. For such large dimensions we would basically measure cache misses rather than performance of the algorithms. Therefore we rewrote all algorithms in single precision. We use the same names `Malcolm`, `LAccu` and `AccSum` as before, so `AccSum` in the following Tables 9.5 and 9.6 refers to `AccSumHugeN`. For increasing dimension fewer and fewer bits can be extracted by `AccSumHugeN` at a time, thus requiring more and more extractions. We tested dimensions from $n = 10,000$ up to $n = 1,280,000$, which is the range of applicability of `AccSumHugeN` where the dimension is too large for `AccSum`. All examples are generated to have condition number $10^8$.

The results for Pentium 4 and Itanium 2 architectures are displayed in Table 9.5, where now the computing time of `SSum`, the ordinary single precision recursive summation, is normed to 1. Moreover, the ratio of computing time for single precision `Sum2` [18] and `XBLAS` [14, 1] to `SSum` is displayed. Note that both deliver a result "as if" calculated in twice single precision (i.e. 48 bits precision). So for condition number $10^8$ we can expect almost full accuracy of the result, but for condition numbers $2^{48} \sim 2 \cdot 10^{14}$ and above we cannot expect a single correct digit of the computed result; hence the comparison to the other algorithms is not quite fair.

As in Part I of this paper we observed a significant performance drop for larger dimensions due to cache misses. So, as also explained in Part I, algorithms `Sum2` and `XBLAS` do not become relatively faster, but the reference `SSum` gets abruptly slower at a certain dimension. As expected, the computing time for `AccSum` grows slowly with increasing dimension. Both `Malcolm` and `LAccu` suffer severely from the small sizes of the internal accumulators. Note that although `AccSum` computes a result of much better quality, namely a faithfully rounded result, it is faster than `XBLAS` up to dimensions where cache misses appear.

The results for AMD Athlon architecture are displayed in Table 9.6, left the ratio of computing times relative

TABLE 9.6
*Test* `AccSumHugeN` *in single precision on AMD Athlon 64 (2.2GHz), GNU gfortran 4.1.1, cond = $10^8$. Left: Measured computing times, time of* `SSumU` *normed to 1, Right: Measured MFlops*

| $n$ | Sum2 | XBLAS | Malcolm | LAccu | AccSum | SSum | SSumU | Malcolm | AccSum |
|---|---|---|---|---|---|---|---|---|---|
| 10,000 | 12.5 | 31.1 | 203.5 | 325.4 | 43.5 | 333 | 2193 | 119 | 1159 |
| 20,000 | 11.0 | 26.1 | 170.0 | 272.3 | 33.8 | 548 | 1838 | 119 | 1253 |
| 40,000 | 9.6 | 23.4 | 151.9 | 243.9 | 37.6 | 551 | 1645 | 119 | 1180 |
| 80,000 | 9.5 | 23.0 | 205.3 | 240.4 | 39.8 | 551 | 1623 | 158 | 1264 |
| 160,000 | 5.3 | 12.6 | 111.4 | 130.1 | 25.9 | 539 | 874 | 157 | 1182 |
| 320,000 | 4.6 | 11.0 | 157.6 | 114.0 | 27.3 | 539 | 767 | 185 | 1207 |
| 640,000 | 4.5 | 10.9 | 156.2 | 112.4 | 34.3 | 536 | 757 | 184 | 1216 |
| 1,280,000 | 4.4 | 10.9 | 222.2 | 112.1 | 45.2 | 538 | 756 | 191 | 1189 |

to recursive summation `SSumU` with unrolled loops, right MFlops. As explained in Part I this is the only architecture out of the three where unrolled loops speed up recursive summation `SSum`. For $n = 10,000$ the speed up is a factor 6.6, i.e. `SSumU` is more than 6 times faster than `SSum`. We programmed also the other algorithms `XBLAS` etc. with unrolled loops, but observed almost no difference.

Again there is a drop in performance at a certain dimension. The MFlop-rates are displayed in the right half of Table 9.6. For `Sum2`, `XBLAS` and `LAccu` those are 1200, 690 and 200, respectively, for all dimensions, so they are not displayed. So the MFlop-rate of `AccSum` is significantly slower than that of `SSumU` for small dimension, and becomes superior for larger dimensions.

Finally we tested `AccSign`. To compute the sign of a sum of floating-point numbers is very simple from the intermediate result of Malcolm's or the long accumulator algorithm. Since a bit representation of the exact sum is available, the sign is immediately available. However, the exact sum is *always* computed, no matter how well- or ill-conditioned the problem is.

For `AccSign` we use the weak "until"-condition with the factor $2^M \text{eps}$ rather than $2^{2M} \text{eps}$ in `AccSum`. We proved that this is weakest possible, just sufficient to guarantee the correct sign. Another competitor is now Priest's doubly compensated summation. The computed approximation `res` of $s := \sum p_i$ satisfies [19, 20] the error estimate $|\text{res} - s| \le 2\text{eps}|s|$, therefore $\text{sign}(\text{res}) = \text{sign}(s)$.

The adapted algorithms by Malcolm and Priest to compute $\text{sign}(s)$ are denoted by `MalcS` and `PriestS`, respectively. To save space we omit results for the long accumulator. The following Table 9.7 shows the results for fixed vector length $n = 1000$. The first column "cond" denotes the condition number of the sum; all computations are performed in double precision. Note again that condition numbers up to $10^{16}$ are the generic case.

As can be seen from Table 9.7, where the time for `AccSum` is normed to 1, the improved stopping criterion in `AccSign` pays for non-extreme condition numbers. The improvement becomes negligible for huge condition numbers. The corresponding ratios for the long accumulator `LAccuS` are at best 38, 20 and 10 for the three architectures and largest condition number $10^{64}$, and for "small" condition number, the generic case, about 110, 51 and 30, respectively.

Finally we tested our algorithms for the special case of zero sums. For Malcolm's, the long accumulator or Priest's algorithm this does not make much difference, but for `AccSign` and `AccSum` it does. In this case the input vector has to be extracted completely, until the final extracted vector is entirely zero. Note that the condition number of a zero sum is infinity.

The results are displayed in Table 9.8, where the time for `AccSign` is normed to 1. All computations are performed in double precision. The first column "Exp. range" depicts $53(1 - \log_{\text{eps}}(e_{max} - e_{min}))$, where

TABLE 9.7
*Measured computing times for* AccSign, $n = 1000$, *for all environments time of* AccSum *normed to 1*

| CPU | Intel Pentium 4 (2.53GHz) | | | Intel Itanium 2 (1.4GHz) | | | AMD Athlon 64 (2.2GHz) | | |
|---|---|---|---|---|---|---|---|---|---|
| Compiler | Intel Visual Fortran 9.1 | | | Intel Fortran 9.0 | | | GNU gfortran 4.1.1 | | |
| cond | PriestS | MalcS | AccSign | PriestS | MalcS | AccSign | PriestS | MalcS | AccSign |
| $10^8$ | 35.4 | 10.1 | 0.55 | 36.8 | 9.5 | 0.71 | 13.8 | 3.6 | 0.43 |
| $10^{16}$ | 35.4 | 10.0 | 0.91 | 36.4 | 9.5 | 0.94 | 13.8 | 3.5 | 0.74 |
| $10^{24}$ | 26.7 | 7.6 | 0.93 | 29.3 | 7.7 | 0.95 | 10.5 | 2.7 | 0.78 |
| $10^{32}$ | 23.6 | 6.8 | 0.95 | 29.4 | 7.8 | 0.95 | 10.5 | 2.7 | 0.78 |
| $10^{40}$ | 18.2 | 5.6 | 0.95 | 18.7 | 5.0 | 0.97 | 7.8 | 2.0 | 0.83 |
| $10^{48}$ | 17.3 | 5.2 | 0.96 | 15.0 | 4.0 | 0.97 | 7.0 | 1.8 | 0.68 |
| $10^{56}$ | 17.4 | 5.3 | 0.96 | 15.0 | 4.0 | 0.97 | 7.0 | 1.8 | 0.83 |
| $10^{64}$ | 17.5 | 5.4 | 0.95 | 15.1 | 4.0 | 0.97 | 5.7 | 1.5 | 0.85 |

TABLE 9.8
*Measured computing times for zero sums,* $n = 1000$, *for all environments time of* AccSign *normed to 1*

| CPU | Intel Pentium 4 (2.53GHz) | | | Intel Itanium 2 (1.4GHz) | | | AMD Athlon 64 (2.2GHz) | | |
|---|---|---|---|---|---|---|---|---|---|
| Compiler | Intel Visual Fortran 9.1 | | | Intel Fortran 9.0 | | | GNU gfortran 4.1.1 | | |
| Exp. range | PriestS | MalcS | LAccu | PriestS | MalcS | LAccu | PriestS | MalcS | LAccu |
| 106 | 33.9 | 9.3 | 63.9 | 24.2 | 7.0 | 32.6 | 15.8 | 3.9 | 22.4 |
| 159 | 32.0 | 9.2 | 61.8 | 24.7 | 7.0 | 32.2 | 15.3 | 3.9 | 22.1 |
| 212 | 21.2 | 6.2 | 42.8 | 16.7 | 4.6 | 21.6 | 10.0 | 2.7 | 14.9 |
| 265 | 14.6 | 4.5 | 30.2 | 12.3 | 3.5 | 16.1 | 7.6 | 2.0 | 11.1 |
| 318 | 12.3 | 3.8 | 24.2 | 9.7 | 2.8 | 12.8 | 6.1 | 1.7 | 8.9 |
| 371 | 10.1 | 3.2 | 20.6 | 8.3 | 2.3 | 10.6 | 4.8 | 1.3 | 7.3 |
| 424 | 8.3 | 2.6 | 17.0 | 7.0 | 2.0 | 9.1 | 4.1 | 1.1 | 6.2 |
| 477 | 7.2 | 2.3 | 15.2 | 6.2 | 1.7 | 8.0 | 3.6 | 1.1 | 5.5 |
| 530 | 6.6 | 2.1 | 13.8 | 5.6 | 1.5 | 7.1 | 3.2 | 1.0 | 4.9 |
| 583 | 6.1 | 1.9 | 12.6 | 4.8 | 1.4 | 6.4 | 2.9 | 0.9 | 4.5 |

$e_{max}$ and $e_{min}$ denote the largest and smallest exponent of the summands, so basically 53 times the number of bits covered by the summands. Zero sums are fortunate for Malcolm's and the long accumulator approach because the *exact* sum *has* to be computed. This is seen from the computational results. Here for AMD Athlon architecture Malcolm's algorithm outperforms AccSign for large exponent range.

Finally we display the achieved MFlop-rates for Itanium and Athlon architecture in Figure 9.1; for Pentium 4 it looks similarly. As can be seen, AccSign achieves for not too large exponent range a little better rate than AccSum. The MFlop-rates for the other algorithms do not change since in any case the exact bit representation is calculated.

**10. Appendix.** Following we prove (6.5), that is we show

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

under the assumptions of Lemma 6.3 and for $\sigma > \frac{1}{2}\texttt{eps}^{-1}\texttt{eta}$. Note that (10.1) implies $R = \text{fl}(\tau_2 - \Delta) = \tau_2 - \Delta$. Since $\sigma$ is a power of 2, we have

$$(10.2) \qquad\qquad \sigma \geq \texttt{eps}^{-1}\texttt{eta} .$$

We already know by (6.4) that

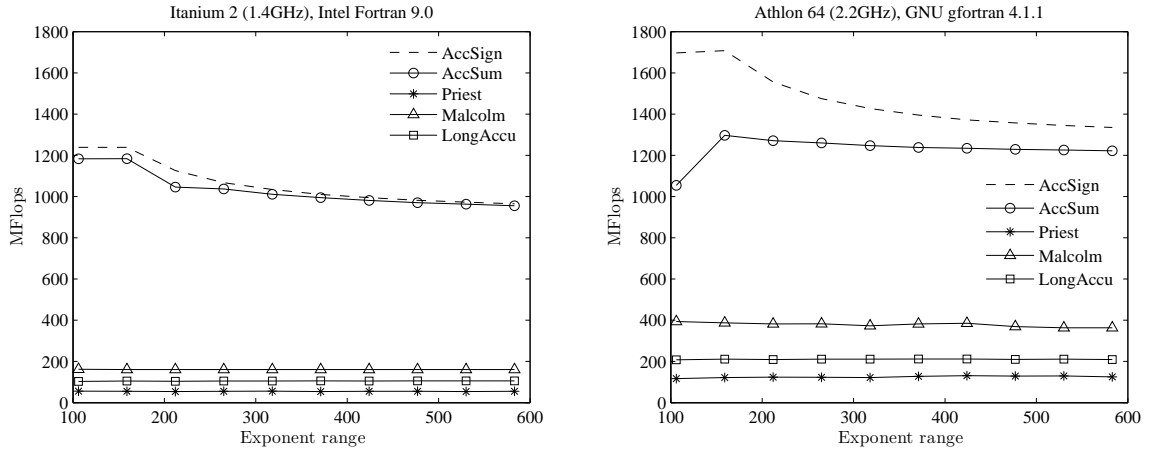$$(10.3) \qquad\qquad \Delta = \texttt{res} - \tau_1 \in \mathbb{F} .$$

FIG. 9.1. *Measured MFlops for varying exponent ranges*

We first prove some facts, namely

$$(10.4) \qquad\qquad \sigma \leq |\tau_1| \quad \Rightarrow \quad \Delta, \tilde{R} \in \mathrm{eps}\sigma\mathbb{Z} \ ,$$

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

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

If $\sigma \leq |\tau_1|$, then $\mathrm{res} = \mathrm{fl}(\tau_1 + \tau_2')$, (2.8) and (2.4) yield $\mathrm{res} \in \mathrm{eps} \cdot \mathrm{ufp}(\tau_1) \subseteq \mathrm{eps}\sigma\mathbb{Z}$, so $\mathrm{res}, \tau_1, \tau_2 \in \mathrm{eps}\sigma\mathbb{Z}$ from $\Delta, \tilde{R} \in \mathrm{eps}\sigma\mathbb{Z}$ and (3.9), proving (10.4).

Furthermore, (3.12), (3.14), (3.9) and the assumption $\frac{1}{2}\mathrm{eps}^{-1}\sigma \leq |\tau_1|$ imply $\sigma \leq 2\mathrm{eps} \cdot \mathrm{ufp}(\tau_1)$ and

$$\begin{aligned} |\tau_2'| &\leq (1+\mathrm{eps})|\tau_2 + \tau_3| \leq |\tau_2| + \mathrm{eps}|\tau_2| + (2^M - 2)\mathrm{eps}\sigma \\ &< |\tau_2| + 2^{M+1}\mathrm{eps}^2\mathrm{ufp}(\tau_1) \leq |\tau_2| + 2^{-M+1}\mathrm{eps} \cdot \mathrm{ufp}(\tau_1) \\ &\leq |\tau_2| + \frac{1}{2}\mathrm{eps} \cdot \mathrm{ufp}(\tau_1) \ , \end{aligned}$$

so that $\tau_1 = \mathrm{fl}(\tau_1 + \tau_2)$ and $\mathrm{res} = \mathrm{fl}(\tau_1 + \tau_2')$ imply (10.5). If $\mathrm{res} \neq \tau_1$, then $\mathrm{sign}(\tau_2) = \mathrm{sign}(\tau_2') = \mathrm{sign}(\mathrm{res} - \tau_1) = \mathrm{sign}(\Delta)$ by (10.3) and the monotonicity of rounding, so (10.6).

PROOF OF (10.1). We use the notation in Lemma 3.5, especially (3.12). We distinguish several cases.

First, assume $|\tau_1| < \sigma$. Then (3.9) yields $\tau_2 \in \mathrm{eps}\sigma\mathbb{Z}$ and $|\tau_2| \leq \mathrm{eps} \cdot \mathrm{ufp}(\tau_1) < \mathrm{eps}\sigma$, so $\tau_2 = 0$ and $\tilde{R} = -\Delta \in \mathbb{F}$ by (10.3).

Second, assume $\sigma \leq |\tau_1| < \frac{3}{5}\mathrm{eps}^{-1}\sigma$. Then (3.12), (2.9), (3.15) and (3.14) imply

$$(10.7) \qquad \begin{aligned} |\tilde{R}| = |\tau_2 - \Delta| &= |\tau_2 - \tau_2' + \delta_1| \leq |\tau_2 - \tau_2'| + \mathrm{eps}|\tau_1 + \tau_2'| \\ &\leq |\tau_3| + \mathrm{eps}|\tau_2 + \tau_3| + \mathrm{eps}|\tau_1| + \mathrm{eps}(1+\mathrm{eps})|\tau_2 + \tau_3| \\ &\leq \mathrm{eps}|\tau_1| + \mathrm{eps}(2+\mathrm{eps})|\tau_2| + (1+\mathrm{eps})^2|\tau_3| \\ &\leq \mathrm{eps}(1+3\mathrm{eps})|\tau_1| + 2^M\mathrm{eps}\sigma \\ &< \sigma \ , \end{aligned}$$

and (2.5), (10.4) and (10.2) prove $\tilde{R} \in \mathbb{F}$.

Henceforth, we may assume without loss of generality $\mathrm{res} \neq \tau_1$ because otherwise $\Delta = \mathrm{res} - \tau_1 = 0$ and $\tilde{R} = \tau_2 \in \mathbb{F}$. For the remaining cases $\frac{3}{5}\mathrm{eps}^{-1}\sigma \leq |\tau_1|$, so (10.5) yields

$$(10.8) \qquad\qquad\qquad \mathrm{res} \in \{\mathrm{pred}(\tau_1), \mathrm{succ}(\tau_1)\} \ .$$

Third, assume $\frac{3}{5}\mathsf{eps}^{-1}\sigma \leq |\tau_1| < \mathsf{eps}^{-1}\sigma$. Then $\mathrm{ufp}(\tau_1) = \frac{1}{2}\mathsf{eps}^{-1}\sigma$ and

$$|\Delta| = |\mathsf{res} - \tau_1| \leq 2\mathsf{eps} \cdot \mathrm{ufp}(\tau_1) = \sigma$$

by Lemma 2.1. But $|\tau_2| \leq \mathsf{eps} \cdot \mathrm{ufp}(\tau_1) = \frac{1}{2}\sigma$, so that $\tau_2 \cdot \Delta \geq 0$ by (10.6) gives

$$|\tilde{R}| = |\tau_2 - \Delta| \leq \sigma \ .$$

Hence (2.5), (10.4) and (10.2) prove $\tilde{R} \in \mathbb{F}$ also for that case.

Fourth and last, assume $|\tau_1| \geq \mathsf{eps}^{-1}\sigma$. We first show that in Algorithm 3.3 (`Transform`) the "repeat-until"-loop is executed only once, i.e. the final value for $m$ is $m = 1$. To establish a contradiction, suppose the final value of $m$ satisfies $m \geq 2$. Then the "until-condition" with $\Phi$ replaced by $2^{2M}\mathsf{eps}$ implies

$$|t^{(m-1)}| < \mathrm{fl}(2^{2M}\mathsf{eps}\sigma_{m-2}) = 2^M\sigma_{m-1} = 2^M\sigma \leq 2^{-M}\mathsf{eps}^{-1}\sigma \leq \frac{1}{4}\mathsf{eps}^{-1}\sigma \ .$$

On the other hand, $\sigma = \sigma_{m-1}$, (3.3) and $|\tau_1| \geq \mathsf{eps}^{-1}\sigma$ imply

$$|t^{(m-1)}| = |\tau_1 + \tau_2 - \tau^{(m)}| > (1 - \mathsf{eps})|\tau_1| - \sigma \geq (1 - \mathsf{eps})\mathsf{eps}^{-1}\sigma - \sigma > \frac{1}{2}\mathsf{eps}^{-1}\sigma \ ,$$

a contradiction. Hence $m = 1$ is the final value of $m$ in Algorithm 3.3 (`Transform`), and (3.3) in Lemma 3.4 yields

(10.9) $$\tau_1 + \tau_2 = \varrho + \tau^{(1)} \quad \text{and} \quad |\tau^{(1)}| \leq (1 - 2^{-M})\sigma \ .$$

Next we prove that (10.8) implies

(10.10) $$|\tau_1| = \mathsf{eps}^{-1}\sigma \quad \text{and} \quad |\Delta| = \sigma \ .$$

Note that by $|\tau_1| \geq \mathsf{eps}^{-1}\sigma$ the distance of $\tau_1$ to its predecessor or successor is at least $\sigma$. We distinguish two cases. First, assume $\tau_1 = \varrho$. Then (10.9) implies $\tau_2 = \tau^{(1)}$ and therefore $|\tau_2| \leq (1 - 2^{-M})\sigma$. Hence (3.12), (3.14) and $2^{2M}\mathsf{eps} \leq 1$ yield

$$|\tau_2'| \leq (1 + \mathsf{eps})|\tau_2 + \tau_3| \leq (1 + \mathsf{eps})\big(1 - 2^{-M} + n\mathsf{eps}\big)\sigma \leq (1 + \mathsf{eps})\big(1 - 2^M\mathsf{eps} + (2^M - 2)\mathsf{eps}\big)\sigma < \sigma \ .$$

The assumption $|\tau_1| \geq \mathsf{eps}^{-1}\sigma$, $\mathsf{res} = \mathrm{fl}(\tau_1 + \tau_2') \neq \tau_1$ and Lemma 2.1 prove $|\tau_1| = \mathsf{eps}^{-1}\sigma$, $|\mathsf{res}| = \mathrm{pred}(|\tau_1|) = |\tau_1| - \sigma$ and $|\Delta| = |\mathsf{res} - \tau_1| = \sigma$, and thus (10.10). Second, assume $\tau_1 \neq \varrho$. In this case $|\tau^{(1)}| < \sigma$ by Lemma 3.4, $|\tau_1| \geq \mathsf{eps}^{-1}\sigma$ and $\tau_1 = \mathrm{fl}(\tau_1 + \tau_2) = \mathrm{fl}(\varrho + \tau^{(1)}) \neq \varrho$ together with Lemma 2.1 imply $|\tau_1| = \mathsf{eps}^{-1}\sigma$, $|\varrho| = \mathrm{pred}(|\tau_1|) = |\tau_1| - \sigma$ and $\tau_1\tau_2 \leq 0$. Now (10.6) yields $\mathrm{sign}(\tau_1) = -\mathrm{sign}(\mathsf{res} - \tau_1)$, and therefore $|\tau_1| > |\mathsf{res}| = |\tau_1| - \sigma$ by (10.8). This proves (10.10).

Finally we have $|\tau_2| \leq \mathsf{eps} \cdot \mathrm{ufp}(\tau_1) = \sigma$, we know $\tau_2\Delta \geq 0$ by (10.6), and together with $|\Delta| = \sigma$ this certifies $|\tilde{R}| = |\tau_2 - \Delta| \leq \sigma$. Therefore (10.4) together with (2.10) yield (10.1). This finishes the proof of (6.5). □

**11. Summary.** We presented algorithms to calculate the sign of a sum, summation with $K$-fold faithfully rounded, with directed rounding and rounded-to-nearest result. The paper contains as well the ingredients to compute a rounded to nearest result in $K$-fold accuracy. All our algorithms use only floating-point addition, subtraction and multiplication in one working precision, no branches in the inner loops and no special operations. Similar algorithms for dot products are easily developed using the error-free transformation `TwoProduct` of a product of two floating-point numbers into a sum (cf. [5], see also [18]). For all algorithms presented in Part I and II of this paper and in [18] we put a Matlab reference code on http://www.ti3.tu-harburg.de/rump .

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.

<div align="center">REFERENCES</div>

[1]  *XBLAS: A Reference Implementation for Extended and Mixed Precision BLAS.* http://crd.lbl.gov/~xiaoye/XBLAS/.

[2]  D. BAILEY, *A Fortran-90 based multiprecision system*, ACM Trans. Math. Software, 21 (1995), pp. 379–387.

[3]  R. BRENT, *A Fortran Multiple-Precision Arithmetic Package*, tech. report, Dept. of Computer Science, Australian National University, Canberra, 1975.

[4]  M. DAUMAS AND D. MATULA, *Validated Roundings of Dot Products by Sticky Accumulation*, IEEE Trans. Comput., 46 (1997), pp. 623–629.

[5]  T. DEKKER, *A Floating-Point Technique for Extending the Available Precision*, Numerische Mathematik, 18 (1971), pp. 224–242.

[6]  L. FOUSSE, G. HANROT, V. LEFÈVRE, P. PÉLISSIER, AND P. ZIMMERMANN, *MPFR: A multiple-precision binary floating-point library with correct rounding*, ACM Trans. Math. Softw., 33 (2007), 15 pages.

[7]  *GNU multiple precision arithmetic library (GMP), version 4.1.2.*, 2003. Code and documentation available at http://swox.com/gmp.

[8]  S. GRAILLAT, P. LANGLOIS, AND N. LOUVET, *Improving the Compensated Horner Scheme with a Fused Multiply and Add*, Proceedings of the ACM symposium on Applied computing, CMS-109, pp. 1323–1327, 2006.

[9]  N. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM Publications, Philadelphia, 2nd ed., 2002.

[10]  *ANSI/IEEE 754-1985, Standard for binary floating-point arithmetic*, 1985.

[11]  D. KNUTH, *The Art of Computer Programming: Seminumerical Algorithms*, vol. 2, Addison Wesley, Reading MA, 1969.

[12]  U. KULISCH AND W. MIRANKER, *Arithmetic Operations in Interval Spaces*, Computing, Suppl., 2 (1980), pp. 51–67.

[13]  P. LANGLOIS AND N. LOUVET, *Solving Triangular Systems more Accurately and Efficiently*, Tech. Report RR2005-02, Laboratoire LP2A, University of Perpignan, 2005.

[14]  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. Software, 28 (2002), pp. 152–205.

[15]  N. LOUVET, *private communication*, 2006.

[16]  M. MALCOLM, *On accurate floating-point summation*, Comm. ACM, 14 (1971), pp. 731–736.

[17]  A. NEUMAIER, *Rundungsfehleranalyse einiger Verfahren zur Summation endlicher Summen*, Zeitschrift für Angew. Math. Mech. (ZAMM), 54 (1974), pp. 39–51.

[18]  T. OGITA, S. RUMP, AND S. OISHI, *Accurate Sum and Dot Product*, SIAM Journal on Scientific Computing (SISC), 26 (2005), pp. 1955–1988.

[19]  D. PRIEST, *Algorithms for Arbitrary Precision Floating Point Arithmetic*, in Proc. 10th Symposium on Computer Arithmetic, P. Kornerup and D. Matula, eds., Grenoble, France, IEEE Computer Society Press, pp. 132–145, 1991.

[20]  ———, *On properties of floating point arithmetics: Numerical stability and the cost of accurate computations*, PhD thesis, Mathematics Department, University of California at Berkeley, CA, USA, 1992. ftp://ftp.icsi.berkeley.edu/pub/theory/priest-thesis.ps.Z.

[21]  S. RUMP, *Kleine Fehlerschranken bei Matrixproblemen*, PhD thesis, Universität Karlsruhe, 1980.

[22]  S. RUMP, T. OGITA, AND S. OISHI, *Accurate Floating-point Summation I: Faithful Rounding.* submitted for publication in SISC, 200–2007.

[23]  S. RUMP AND P. ZIMMERMANN, *Interval operations in rounding to nearest.* submitted for publication, 2007.

[24]  J. SHEWCHUK, *Adaptive precision floating-point arithmetic and fast robust geometric predicates*, Discrete Comput. Geom., 18 (1997), pp. 305–363.

[25]  Y. ZHU AND W. HAYES, *Fast, Guaranteed-Accurate Sums of many Floating-Point Numbers*, in Proc. of the RNC7 Conference on Real Numbers and Computers, G. Hanrot and P. Zimmermann, eds., 2006, pp. 11–22.

[26]  Y. ZHU, J. YONG, AND G. ZHENG, *A New Distillation Algorithm for Floating-Point Summation*, SIAM J. Sci. Comput. (SISC), 26 (2005), pp. 2066–2078.

[27]  G. ZIELKE AND V. DRYGALLA, *Genaue Lösung linearer Gleichungssysteme*, GAMM Mitt. Ges. Angew. Math. Mech., (2003), pp. 7–108.