



HAL
open science

A Class of Fast and Accurate Summation Algorithms

Pierre Blanchard, Nicholas J Higham, Théo Mary

► **To cite this version:**

Pierre Blanchard, Nicholas J Higham, Théo Mary. A Class of Fast and Accurate Summation Algorithms. SIAM Journal on Scientific Computing, In press. hal-02491081

HAL Id: hal-02491081

<https://hal.science/hal-02491081>

Submitted on 25 Feb 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A CLASS OF FAST AND ACCURATE SUMMATION ALGORITHMS*

PIERRE BLANCHARD[†], NICHOLAS J. HIGHAM[‡], AND THEO MARY[§]

Abstract. The need to sum floating-point numbers is ubiquitous in scientific computing. Standard recursive summation of n summands, often implemented in a blocked form, has a backward error bound proportional to nu , where u is the unit roundoff. With the growing interest in low precision floating-point arithmetic and ever increasing n in applications, computed sums are more likely to have insufficient accuracy. We propose a class of summation algorithms called FABsum (for “fast and accurate block summation”) that applies a fast summation algorithm (such as recursive summation) blockwise and then sums the partial sums using an accurate summation algorithm (such as compensated summation, or recursive summation in higher precision). We give a rounding error analysis to show that FABsum with a fixed block size b has a backward error bound $(b+1)u + O(u^2)$, which is independent of n to first order. Our computational experiments show that with a suitable choice of b (independent of n) FABsum can deliver substantially more accurate results than blocked recursive summation, with only a modest drop in performance. FABsum is especially attractive for low precisions, where it can provide correct digits for much larger n than recursive summation.

Key words. summation, recursive summation, compensated summation, blocked summation, rounding error analysis, floating-point arithmetic, numerical linear algebra.

AMS subject classifications. 65G50, 65F05

1. Introduction. Summation is a key computational task at the heart of many numerical algorithms, most notably numerical linear algebra kernels involving inner products, such as matrix–vector or matrix–matrix multiplications, matrix factorizations, and the solution of linear systems. In floating-point arithmetic, summation incurs rounding errors. Backward error bounds for the basic summation algorithms are proportional to nu , where n is the number of summands and u the unit roundoff, and thus they grow linearly with n .

While compensated summation algorithms achieving backward error bounds that do not grow (or grow slowly) with n are known [10, sec. 4.3], they are computationally expensive. Indeed running a basic summation algorithm in higher precision (e.g., using double precision rather than single precision) may provide a better performance/accuracy tradeoff. Compensated summation algorithms are thus only used when the highest precision available is not enough for the application at hand, and they are generally not available in optimized libraries such as the Intel Math Kernel Library or the NAG Library.

The rise of large-scale, low-precision computations presents new challenges. On the one hand, modern supercomputing allows larger and larger problems to be solved, and we must routinely evaluate sums of 10^8 or more terms. On the other hand, the use of low floating-point precisions (such as IEEE half precision, for which $u \approx 4.88 \times 10^{-4}$) is increasingly common [2], [8], [12]. For such large sizes and low precisions, error bounds of order nu exceed 1, and so do not provide any meaningful information. This

*Version of February 19, 2020. **Funding:** This work was supported by Engineering and Physical Sciences Research Council grant EP/P020720/1, The MathWorks, and the Royal Society. The opinions and views expressed in this publication are those of the authors, and not necessarily those of the funding bodies.

[†]Department of Mathematics, The University of Manchester, Manchester M13 9PL, UK (pierre.blanchard00@gmail.com)

[‡]Department of Mathematics, The University of Manchester, Manchester, M13 9PL, UK (nick.higham@manchester.ac.uk, <http://www.maths.manchester.ac.uk/~higham>)

[§]Department of Mathematics, The University of Manchester, Manchester M13 9PL, UK (theo.mary@manchester.ac.uk)

creates a dilemma, as we need to choose between benefiting from the speed, lower energy requirements, and reduced data movement of low precision arithmetic and being able to accurately solve large problems.

In this work we tackle this dilemma by introducing a new class of summation algorithms that excels in both performance and accuracy. These algorithms achieve backward error bounds that do not grow with n (that is, the bounds are of the form $cu + O(u^2)$, where c is a moderate constant independent of n) and, at the same time, they can deliver a similar performance to optimized, standard algorithms by performing an arbitrarily low number of extra flops (e.g., less than a 1% overhead) and by allowing efficient implementation on modern computers.

We first review existing summation algorithms in section 2. Then, in section 3, we present the new class of algorithms, called **FABsum**, perform a rounding error analysis for algorithms in this class, and analyze their cost in a general framework. In section 4 we apply **FABsum** to several key numerical linear algebra algorithms. In section 5, we assess the practical performance and accuracy of **FABsum** through numerical experiments. We provide our concluding remarks in section 6.

2. Existing summation algorithms. We begin by briefly reviewing existing algorithms to compute the sum $s = \sum_{i=1}^n x_i$. We denote by \hat{s} the computed value of s in floating-point arithmetic. We use the standard model of floating-point arithmetic [10, sec. 2.2]

$$(2.1) \quad \text{fl}(a \text{ op } b) = (a \text{ op } b)(1 + \delta), \quad |\delta| \leq u, \quad \text{op} \in \{+, -, \times, /\}.$$

Note that this model does not account for the possibility of underflow (or overflow). However, an addition or subtraction incurring underflow is necessarily exact [9], [10, Prob. 2.19] and therefore our results for summation hold even in presence of underflow. This does not extend, though, to the more general algorithms based on summation in section 4.

We note that sharper error bounds can be obtained in some cases by making the stronger assumption that the floating-point arithmetic uses round to nearest [16], [19], [24]. However, we prefer to restrict to the model (2.1) in order that our results have the widest possible applicability, for example to directed or stochastic rounding.

We mention that summation algorithms have been proposed that achieve reproducibility (see, e.g., [4], [5]). Reproducibility and accuracy are two distinct objectives, since even very accurate algorithms can yield different results on different runs. In this article, we target accuracy rather than reproducibility.

2.1. Recursive summation. The most basic algorithm to compute s is recursive summation, which starts with $s = x_1$ and computes

$$s \leftarrow s + x_i, \quad i = 2 : n.$$

The computed sum satisfies [10, sec. 4.2]

$$(2.2) \quad \hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq \gamma_{n-1} = (n-1)u + O(u^2),$$

where $\gamma_n = nu/(1 - nu)$ for $nu < 1$. This backward error bound grows linearly with n and is almost attainable [10, Prob. 4.2].

2.2. Blocked summation. In order to allow a sum to be computed in parallel, high-performance libraries divide the summands into n/b blocks of size b . From now on, n is assumed to be a multiple of b for simplicity. First, the local sums

$$(2.3) \quad s_i = \sum_{j=(i-1)b+1}^{ib} x_j, \quad i = 1 : n/b,$$

are computed with recursive summation, and then $s = \sum_{i=1}^{n/b} s_i$ is also obtained with recursive summation.

By applying (2.2) to s_i and s we find that the computed \hat{s} from blocked summation satisfies

$$(2.4) \quad \hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq (b + n/b - 2)u + O(u^2),$$

and so the backward error bound depends on the block size b . The constant $b + n/b - 2$ in the bound is minimized for $b = \sqrt{n}$, for which it equals $2\sqrt{n} - 2$. However, for performance reasons b is often chosen to be a moderate constant (such as 128 or 256), in which case the error bound for blocked summation still grows linearly with n .

2.3. Pairwise summation. Pairwise summation (also known as tree summation or binary cascade summation) generalizes blocked summation by recursively computing the local sums (2.3) with blocked summation, continuing until sums of two terms are obtained. There are thus at most $\lceil \log_2 n \rceil$ levels of recursion, where $\lceil \cdot \rceil$ denotes the ceiling function, which leads to the backward error result [10, sec. 4.2]

$$(2.5) \quad \hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq \lceil \log_2 n \rceil u + O(u^2),$$

in which the bound grows with n at a much slower rate than for blocked summation. However, pairwise summation is not well suited for an efficient implementation on parallel computers, and it usually delivers poor performance compared with blocked summation [3].

2.4. Compensated summation. The principle of compensation is based on the fact that for any pair (a, b) of floating-point numbers and their computed sum $\hat{s} = fl(a + b)$, there exists a floating-point number e such that

$$(2.6) \quad a + b = \hat{s} + e.$$

In other words, the error in a floating-point addition is itself a floating-point number.

Compensated summation algorithms approximate the error term e at each step of recursive summation in order to correct the computed sum. Many such algorithms have been designed and analyzed over the years. Here we focus on Algorithm 2.1 [17], [22], and we note that the augmented arithmetic operations in the 2019 revision of the IEEE 754 standard for floating-point arithmetic [15] can be exploited in the algorithm. The backward error result for this algorithm is [7, Thm. 8], [18, Ex. 19, sec. 4.2.2]

$$(2.7) \quad \hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq 2u + O(u^2),$$

Algorithm 2.1 Compensated summation.

This algorithm takes as input n summands x_i and computes their sum by compensated summation.

```
1:  $s = 0, e = 0$ 
2: for  $i = 1 : n$  do
3:    $z = s$ 
4:    $y = x_i + e$ 
5:    $s = z + y$ 
6:    $e = (z - s) + y$ 
7: end for
```

which is an almost ideal result.

Algorithm 2.1 requires 3 extra flops per loop iteration, which is typically more expensive than simply switching to a higher precision. For this reason, compensated summation is usually worth considering only when computations are already taking place at the highest precision available on the target hardware.

3. Fast and accurate blocked summation. Assume that we have at our disposal two summation algorithms: one that is fast, referred to as the **FastSum** algorithm, and one that is accurate, referred to as the **AccurateSum** algorithm. In Algorithm 3.1 we describe **FABsum**, a new blocked summation algorithm that exploits **FastSum** and **AccurateSum**. The key idea is to use **FastSum** to compute the local sums s_i of b numbers and **AccurateSum** to sum the local sums. This idea is motivated by the observations that the bulk of the computation is in computing the local sums if the block size is chosen appropriately and that rounding error growth can be attenuated by accumulating the local sums more accurately. We can therefore expect that the new algorithm will be almost as fast as **FastSum** and almost as accurate as **AccurateSum**. The rest of this section is devoted to proving this expectation by analyzing the rounding errors and the cost of Algorithm 3.1 and illustrating these analyses with some practical choices of **FastSum** and **AccurateSum**.

Algorithm 3.1 (**FABsum**) Fast and accurate blocked summation algorithms.

This algorithm takes as input n summands x_i , a block size b , and two summation algorithms **FastSum** and **AccurateSum**. It returns the sum $s = \sum_{i=1}^n x_i$.

```
1: for  $i = 1 : n/b$  do
2:   Compute  $s_i = \sum_{j=(i-1)b+1}^{ib} x_j$  with FastSum.
3: end for
4: Compute  $s = \sum_{i=1}^{n/b} s_i$  with AccurateSum.
```

We note the special cases $b = 1$, which is entirely **AccurateSum**, and $b = n$, which is entirely **FastSum**. We are interested in values of b between these two extremes: b should be large enough to give speed benefits but small enough to give error bounds independent of n to first order. In the next two sections we give error and cost analyses that guide the choice of b .

Our general framework includes some previously proposed algorithms. For example, taking **AccurateSum** to be pairwise summation leads to the superblock algorithm from [3]. Similarly, [1] presents a “selective compensation” algorithm that amounts to using compensated summation as **AccurateSum**. Our analysis below generalizes the

analyses of these papers, considers other possible `AccurateSum` choices, and extends the analysis to numerical linear algebra kernels in section 4.

3.1. Rounding error analysis. To carry out a rounding error analysis of Algorithm 3.1 we need to make some assumptions on the error bounds associated with algorithms `FastSum` and `AccurateSum`. For any sum of the form $s = \sum_{i=1}^n x_i$, we assume that the computed \hat{s} from `FastSum` satisfies

$$(3.1) \quad \hat{s} = \sum_{i=1}^n x_i(1 + \mu_i^f), \quad |\mu_i^f| \leq \varepsilon_f(n)$$

and the computed \hat{s} from `AccurateSum` satisfies

$$(3.2) \quad \hat{s} = \sum_{i=1}^n x_i(1 + \mu_i^a), \quad |\mu_i^a| \leq \varepsilon_a(n),$$

where $\varepsilon_f(n)$ and $\varepsilon_a(n)$ are $O(u)$ and depend on n and u . With these assumptions, we have the following backward error result.

THEOREM 3.1. *Let $s = \sum_{i=1}^n x_i$ be computed by Algorithm 3.1. The computed \hat{s} satisfies*

$$\hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq \varepsilon(n, b) = \varepsilon_f(b) + \varepsilon_a(n/b) + \varepsilon_f(b)\varepsilon_a(n/b).$$

Proof. Each of the local sums s_i computed at line 2 of Algorithm 3.1 satisfies

$$(3.3) \quad \hat{s}_i = \sum_{j=(i-1)b+1}^{ib} x_j(1 + \mu_j^f), \quad |\mu_j^f| \leq \varepsilon_f(b).$$

Let $t = \sum_{i=1}^{n/b} \hat{s}_i$ be the sum of the computed local sums. Then the computed \hat{t} , which is the overall computed sum \hat{s} , satisfies

$$(3.4) \quad \hat{t} = \hat{s} = \sum_{i=1}^{n/b} \hat{s}_i(1 + \mu_i^a), \quad |\mu_i^a| \leq \varepsilon_a(n/b).$$

Combining (3.3) and (3.4) and defining $\tilde{\mu}_i^a = \mu_{\lceil i/b \rceil}^a$, we obtain

$$\hat{s} = \sum_{i=1}^n x_i(1 + \mu_i^f)(1 + \tilde{\mu}_i^a), \quad |\mu_i^f| \leq \varepsilon_f(b), \quad |\tilde{\mu}_i^a| \leq \varepsilon_a(n/b).$$

Defining $\mu_i = \mu_i^f + \tilde{\mu}_i^a + \mu_i^f \tilde{\mu}_i^a$, the result follows. \square

We can interpret Theorem 3.1 as follows. If ε_a is of order u^2 (such as when `AccurateSum` uses doubled precision, that is, with a unit roundoff u^2 corresponding to a significand twice as large), then $\varepsilon(n, b) = \varepsilon_f(b) + O(u^2)$ is independent of n to first order and thus does not grow with n . In fact, even if ε_a is of order u , as long as it does not grow with n to first order (such as when `AccurateSum` uses compensation), $\varepsilon(n, b)$ does not grow with n to first order either.

To explore this point more precisely we provide an expression for $\varepsilon(n, b)$ in the following cases. Assume `FastSum` corresponds to recursive summation, so that $\varepsilon_f(b) = (b-1)u + O(u^2)$.

- If `AccurateSum` uses recursive summation in doubled precision with a final rounding back to the working precision, then we have $\varepsilon_a(n/b) = u + O(u^2)$ and an overall error bound

$$(3.5) \quad \varepsilon(n, b) = bu + O(u^2).$$

- If `AccurateSum` uses compensated summation, we have $\varepsilon_a(n/b) = 2u + O(u^2)$ by (2.7) and an overall bound

$$(3.6) \quad \varepsilon(n, b) = (b + 1)u + O(u^2).$$

- If `AccurateSum` uses pairwise summation, we have $\varepsilon_a(n/b) = \lceil \log_2(n/b) \rceil u + O(u^2)$ by (2.5) and an overall bound

$$(3.7) \quad \varepsilon(n, b) = (b - 1 + \lceil \log_2(n/b) \rceil)u + O(u^2).$$

In all cases we thus obtain an overall error bound that does not grow (or, in the case of pairwise summation, grows very slowly) with n .

3.1.1. Accuracy versus stability. Accurate summation means achieving a small forward error, but our error analysis bounds the backward error. We will show that these concepts are closely related for summation. The backward error of an approximation \hat{s} to $s = \sum_{i=1}^n x_i$ is

$$(3.8) \quad \eta(\hat{s}) = \min \left\{ \max |\delta_i| : \hat{s} = \sum_{i=1}^n x_i(1 + \delta_i) \right\}.$$

A formula for the backward error can be obtained as a special case of the Oettli–Prager theorem [10, Thm. 7.3], [23]:

$$(3.9) \quad \eta(\hat{s}) = \frac{|\hat{s} - \sum_{i=1}^n x_i|}{\sum_{i=1}^n |x_i|}.$$

The numerator is the absolute forward error in the computed sum. The relative forward error is a factor $\sum_{i=1}^n |x_i| / \sum_{i=1}^n x_i$ larger than $\eta(\hat{s})$ —a factor that depends on the data but not the algorithm. This factor is precisely the condition number

$$\text{cond}(s) = \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{|\sum_{i=1}^n (x_i + \Delta x_i) - \sum_{i=1}^n x_i|}{\varepsilon |\sum_{i=1}^n x_i|} : |\Delta x_i| \leq \varepsilon |x_i|, i = 1 : n \right\}.$$

3.1.2. On second order contributions. Up to now we have considered expansions of the error bounds up to first order in u . However, when low precision arithmetic is being used it is likely that the $O(u^2)$ terms will become significant even for moderate n .

If the arithmetic uses round to nearest we can drop these second order terms for recursive summation, as proved in [16], [19], [24]. However, as noted in section 2, we wish to keep the analysis more general. Using the general model (2.1) we now obtain the second order terms explicitly in order to assess their potential influence on the error bounds of `FABsum`.

Let us start with the error bounds given in section 2. In the case of recursive summation, if $(n - 1)u < 1$ a second order bound is found by using the expansion of γ_{n-1} at $u = 0$ in (2.2):

$$\hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq (n - 1)u + (n - 1)^2 u^2 + O(u^3).$$

Here, the second order term becomes significant when $(n-1)u$ approaches 1, that is, when the bound becomes invalid and hence of no use. Similarly, for recursive summation in doubled precision (with unit roundoff u^2) with a final rounding to working precision u , we have

$$\hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq u + (n-1)u^2 + O(u^3).$$

Finally, an explicit second order backward error bound for compensated summation can be found in [7, Thm. 8]:

$$\hat{s} = \sum_{i=1}^n x_i(1 + \mu_i), \quad |\mu_i| \leq 2u + 2(2n+1)u^2 + O(u^3).$$

By Theorem 3.1, the second order terms of `FABsum` can be obtained by adding the second order terms of $\varepsilon_f(b)$ and $\varepsilon_a(n/b)$ to the product of their first order terms. We assume that `FastSum` uses recursive summation at the working precision, so that $\varepsilon_f(b) = (b-1)u + (b-1)^2u^2 + O(u^3)$.

- If `AccurateSum` uses recursive summation in doubled precision, we have $\varepsilon_a(n/b) = u + (n/b-1)u^2 + O(u^3)$ and an overall error bound

$$\varepsilon(n, b) = bu + [n/b - 1 + (b-1)^2 + (b-1)]u^2 + O(u^3).$$

- If `AccurateSum` uses compensated summation, we have $\varepsilon_a(n/b) = 2u + 2(2n/b+1)u^2 + O(u^3)$ and an overall bound

$$\varepsilon(n, b) = (b+1)u + [4n/b + 2 + (b-1)^2 + 2(b-1)]u^2 + O(u^3).$$

We will not derive a higher-order bound for pairwise summation as the bound (3.7) depends on n to first order already. For large n , n/b is in practice much larger than b and so we have, roughly,

$$(3.10) \quad \varepsilon(n, b) = bu + (n/b)u^2 + O(u^3)$$

for doubled precision and

$$(3.11) \quad \varepsilon(n, b) = (b+1)u + (4n/b)u^2 + O(u^3)$$

for compensated summation. The order u^2 terms become significant when n exceeds the critical values b^2/u for doubled precision and $b^2/(4u)$ for compensated summation. For instance, in IEEE half precision arithmetic ($u = 2^{-11} \approx 4.88 \times 10^{-4}$) and for a block size $b = 128$, the second order terms in the bounds are significant for n larger than 3.3×10^7 with doubled precision and 8.3×10^6 with compensation.

3.2. Cost analysis. We now analyze the cost $C(n, b)$ in flops of Algorithm 3.1. Let C_f and C_a denote the costs of algorithms `FastSum` and `AccurateSum`, respectively. Then

$$C(n, b) = \frac{n}{b}C_f(b) + C_a\left(\frac{n}{b}\right).$$

In particular, if the costs C_f and C_a are linear functions of the number of summands, as is often the case, $C(n, b)$ simplifies to

$$C(n, b) = C_f(n) + \frac{1}{b}C_a(n) + O\left(\frac{n}{b}\right).$$

Therefore the cost of Algorithm 3.1 can be made arbitrarily close to that of `FastSum` by increasing the block size b .

For example, assume `FastSum` is recursive summation, whose cost is $C_f(n) = n - 1$, and `AccurateSum` is compensated summation, whose cost is $C_a(n) = 4n - 2$. Then $C(n, b)$ is equal to $n - 1 + (4n - 2)/b$, that is, it requires only $(4n - 2)/b$ extra flops compared with recursive summation, which is a very small overhead for typical choices of block sizes. For instance, for $b = 400$, this represents an overhead of only 1%. Similar results hold when `AccurateSum` uses recursive summation with doubled precision (the use of doubled precision increases the cost of recursive summation by a constant factor, usually 2) or pairwise summation (which does not change the number of required flops, but may decrease the speed due to poor efficiency on parallel computers).

As long as the parameter b is both reasonably large (to keep the overhead cost limited) and constant (to avoid error growth), we therefore have the freedom to tune it in order to achieve the highest possible performance on a given target architecture, taking into account the cache size, for example.

4. Application to numerical linear algebra. Now we apply our new `FABsum` algorithm (Algorithm 3.1) within some important summation-based kernels in numerical linear algebra, namely inner products, matrix–vector and matrix–matrix products, matrix factorizations, and the solution of linear systems. The core computation is an inner product, and to be precise we specify in Algorithm 4.1 how `FABsum` is used.

Algorithm 4.1 Inner product via `FABsum`.

This algorithm takes as input vectors $x, y \in \mathbb{R}^n$ and a block size b and returns the inner product $x^T y$.

- 1: Apply Algorithm 3.1 to the summands $x_i y_i$, $i = 1 : n$.
-

We now derive rounding error bounds for the resulting algorithms. Recall that $\varepsilon_f(n)$ and $\varepsilon_a(n)$ are the error constants in (3.1) and (3.2) for `FastSum` and `AccurateSum`, respectively.

THEOREM 4.1 (inner products). *Let $x, y \in \mathbb{R}^n$ and let $s = x^T y$ be computed by Algorithm 4.1. The computed \hat{s} satisfies*

$$\hat{s} = \sum_{i=1}^n x_i y_i (1 + \mu_i), \quad |\mu_i| \leq u + \varepsilon_f(b) + \varepsilon_a(n/b) + O(u^2).$$

Proof. The initial products $z_i = x_i y_i$ each incur one rounding error. The rest of the proof follows from the application of Theorem 3.1 to the computation of $\sum_{i=1}^n z_i$. \square

In the next result we consider matrix–vector and matrix–matrix products computed by the usual inner product formulas. Here, $|A|$ denotes the matrix $(|a_{ij}|)$ and inequalities involving matrices or vectors are to be understood componentwise.

THEOREM 4.2 (matrix–vector and matrix–matrix products). *Let $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$, and $x \in \mathbb{R}^n$. If $y = Ax$ is computed via inner products using Algorithm 4.1 then*

$$\hat{y} = (A + \Delta A)x, \quad |\Delta A| \leq (u + \varepsilon_f(b) + \varepsilon_a(n/b) + O(u^2)) |A|.$$

If $C = AB$ is computed via inner products using Algorithm 4.1 then

$$|\hat{C} - C| \leq (u + \varepsilon_f(b) + \varepsilon_a(n/b) + O(u^2)) |A| |B|.$$

Proof. The proof follows directly from Theorem 4.1. \square

In order to analyze the solution of linear systems we need the next lemma. For convenience in stating the bounds we will assume that $\varepsilon_f(n) \geq u$ and $\varepsilon_a(n) \geq u$, which is true for all the methods under consideration here.

LEMMA 4.3. *If $y = (c - \sum_{i=1}^{k-1} a_i b_i)/b_k$ is computed by applying Algorithm 4.1 to the summation term then the computed \hat{y} satisfies*

$$(4.1) \quad b_k \hat{y}(1 + \mu_0) = c - \sum_{i=1}^{k-1} a_i b_i (1 + \mu_i),$$

where $|\mu_i| \leq u + \varepsilon_f(b) + \varepsilon_a((k-1)/b) + O(u^2)$ for all i .

Proof. By Theorem 3.1 and (2.1) we have

$$\hat{y} = b_k^{-1} \left(c - \sum_{i=1}^{k-1} a_i b_i (1 + \delta_i)(1 + \mu'_i) \right) (1 + \delta_k)(1 + \delta_{k+1}),$$

where

$$|\delta_i| \leq u, \quad i = 1:k+1, \quad |\mu'_i| \leq \varepsilon_f(b) + \varepsilon_a((k-1)/b) + \varepsilon_f(b)\varepsilon_a((k-1)/b).$$

Therefore (4.1) holds with $1 + \mu_0 = ((1 + \delta_k)(1 + \delta_{k+1}))^{-1}$, which implies $|\mu_0| \leq 2u + O(u^2)$, and

$$1 + \mu_i = 1 + \delta_i + \mu'_i + \delta_i \mu'_i,$$

which implies

$$|\mu_i| \leq u + \varepsilon_f(b) + \varepsilon_a((k-1)/b) + O(u^2). \quad \square$$

We can now derive results for the solution of linear systems by LU factorization. The next result is immediate from Lemma 4.3.

THEOREM 4.4. *Let the triangular system $Tx = b$, where $T \in \mathbb{R}^{n \times n}$, be solved by substitution using Algorithm 4.1. Then the computed solution \hat{x} satisfies*

$$(T + \Delta T)\hat{x} = b, \quad |\Delta T| \leq (u + \varepsilon_f(b) + \varepsilon_a((n-1)/b) + O(u^2))|T|.$$

For the next two results we assume that the Doolittle form of LU factorization is used (see, e.g., [10, Alg. 9.2]) and that Algorithm 4.1 is used in the inner products therein.

THEOREM 4.5. *If LU factorization applied to $A \in \mathbb{R}^{n \times n}$ runs to completion then the computed LU factors \hat{L} and \hat{U} satisfy*

$$A + \Delta A = \hat{L}\hat{U}, \quad |\Delta A| \leq (u + \varepsilon_f(b) + \varepsilon_a(n/b) + O(u^2))|\hat{L}||\hat{U}|.$$

Proof. The result is obtained by applying Lemma 4.3 to the equations that determine L and U , and is analogous to the proof of [10, Thm. 9.3]. \square

THEOREM 4.6. *Let $A \in \mathbb{R}^{n \times n}$ and suppose LU factorization produces computed LU factors \hat{L} , \hat{U} , and a computed solution \hat{x} to $Ax = b$. Then*

$$(4.2) \quad (A + \Delta A)\hat{x} = b, \quad |\Delta A| \leq 3(u + \varepsilon_f(b) + \varepsilon_a(n/b) + O(u^2))|\hat{L}||\hat{U}|.$$

Proof. The result is obtained by combining Theorems 4.4 and 4.5, and is analogous to the proof of [10, Thm. 9.4]. \square

When recursive summation is used (which corresponds to $b = n$ in Algorithm 3.1, with `FastSum` equal to recursive summation), we have $\varepsilon_a = 0$ and $\varepsilon_f(b) = \varepsilon_f(n) = (n - 1)u + O(u^2)$, so the error constants in Theorems 4.1–4.2 and 4.4–4.6 have first order terms proportional to n . But as shown in section 3.1, if `FABsum` is used with `FastSum` equal to recursive summation and `AccurateSum` equal to either compensated summation or recursive summation in doubled precision then the first order terms in Theorems 4.1–4.2 and 4.4–4.6 are of order bu and hence independent of n . `FABsum` therefore offers error constants roughly n/b times smaller than those for recursive summation.

5. Numerical experiments. We have carried out an extensive set of numerical experiments to assess the performance and accuracy of our new `FABsum` algorithm (Algorithm 3.1) with different `AccurateSum` choices and to compare it with existing algorithms.

In section 5.1 we experiment with summation using MATLAB R2018b. In section 5.2 we present performance and accuracy results with PLASMA [6], a state-of-the-art numerical linear algebra library that we have modified by integrating `FABsum`.

In all the experiments, the working precision is single ($u = 2^{-24} \approx 5.96 \times 10^{-8}$) and the “exact” quantities appearing in the error bound formulas for inner products and matrix–matrix products are computed in double precision.

5.1. Summation. We first experiment with summation with MATLAB. We compute $s = \sum_{i=1}^n x_i$, where x is a randomly generated vector. While the behavior of the error as a function of n does not strongly depend on the type of distribution (such as random uniform or normal), it does depend on the mean of the entries x_i [13]. For this reason, we compare the two distributions

- random uniform $[0, 1]$: `x = rand(m,n)`,
- random uniform $[-1, 1]$: `x = (rand(m,n)-0.5)*2`.

In order to make the experiments reproducible, we seed the random number generator at the beginning of the scripts generating each figure of this section. We have made these scripts available online¹. For each problem size n , we run the same experiment $N_{test} = 10$ times and plot the value of the maximum backward error (3.8).

Results using single precision and a block size $b = 128$ are displayed in Figure 5.1, where we compare blocked summation (one of the fastest summation algorithms), compensated summation (Algorithm 2.1, one of the most accurate summation algorithms), and `FABsum` (Algorithm 3.1), which in this case uses recursive summation for `FastSum` and compensated summation for `AccurateSum`. Clearly, `FABsum` is much more stable than blocked summation, by up to three orders of magnitude for large n . It yields a backward error that does not grow with n and is almost ideally small for $[0, 1]$ vectors, matching the compensated summation error. For $[-1, 1]$ vectors, compensated summation remains more accurate, but only by a constant factor always less than 10.

The backward error bound for `FABsum` is $(b + 1)u + O(u^2)$ (see (3.6)) versus $2u + O(u^2)$ for compensated summation, but we see a ratio in the actual backward errors roughly $\sqrt{b}/2$ rather than $b/2$. This is not surprising because the probabilistic error analysis in [11] shows that if the rounding errors are independent random variables of zero mean then an error bound of order the square root of the worst-case bound holds with high probability.

The next experiment uses IEEE half precision (fp16, $u = 2^{-11} \approx 4.88 \times 10^{-4}$),

¹<https://gitlab.com/nla-grp/FABsum>.

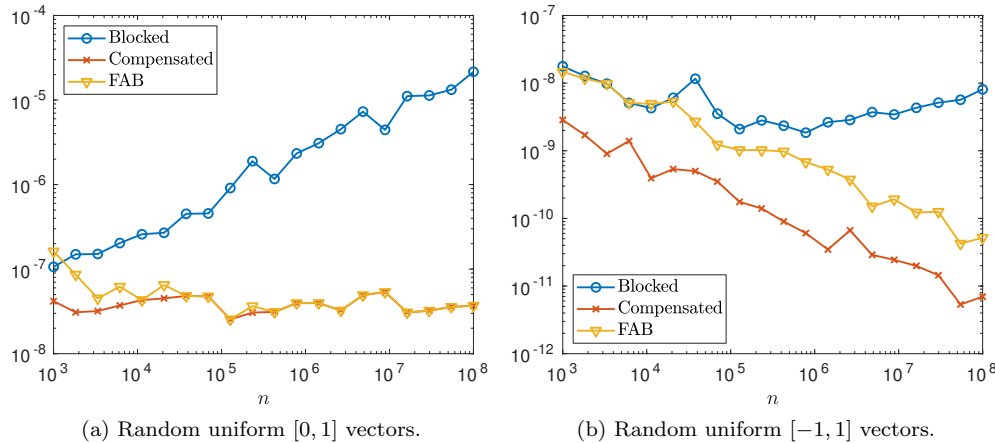


Fig. 5.1: Backward error for the sum $\sum_{i=1}^n x_i$ computed in single precision, for a vector x with random uniform entries and block size $b = 128$.

bfloat16 ($u = 2^{-8} \approx 3.91 \times 10^{-3}$), and quarter precision ($u = 2^{-4} \approx 6.25 \times 10^{-2}$)—the latter is not standard, but this choice of precision for 8-bit words is suggested in [21]. These lower precisions are simulated with the rounding function `chop.m` from [14]. Since our objective is to assess the effect of rounding errors, in order to avoid overflow interfering with the results we only simulate the significand of these low precision formats; for the exponent, we keep the same number of bits as double precision. Figure 5.2 compares blocked summation with FABsum with recursive summation for `FastSum` and two choices for `AccurateSum`: compensated summation and recursive summation in extended precision $u_e = 2^{-24}$ (that is, single precision). As can be seen in the figure, for such low precisions blocked summation leads to a backward error that quickly reaches 1, at which point the computed sum is meaningless from a backward error point of view. Since in this experiment x is a nonnegative vector, one important difficulty is that the sum s keeps increasing as it is computed. For large n , the sum may become so large that the following increments do not increase its value in floating-point arithmetic. This phenomenon, which we call *stagnation*, leads to a dramatic increase of the error. The reason is that the rounding errors incurred during stagnation are all negative, which makes the worst-case error bound close to being sharp [11, Sec. 4.2.1] and also invalidates the model upon which the probabilistic analysis in [11] is based. Since the entries of x are bounded by 1, we know that s (in line 4 of Algorithm 3.1) must necessarily stagnate after its value exceeds b/u ; note that stagnation can, however, occur for smaller values of s when the increments x_i are small. For random uniform summands in $[0, 1]$, the value of s should be about $n/2$. This explains the dramatic increase of the error for blocked summation observed in Figure 5.2 when n becomes larger than $2b/u$, that is, 131072, 16384, and 1024 in fp16, bfloat16, and fp8 arithmetic, respectively.

We now discuss the behavior of FABsum with respect to stagnation for different choices of `AccurateSum`. For fp8 arithmetic, stagnation occurs for FABsum with compensated summation, eventually leading to a sharp increase of the error. However, when extended (single) precision is used, s does not stagnate because it never exceeds

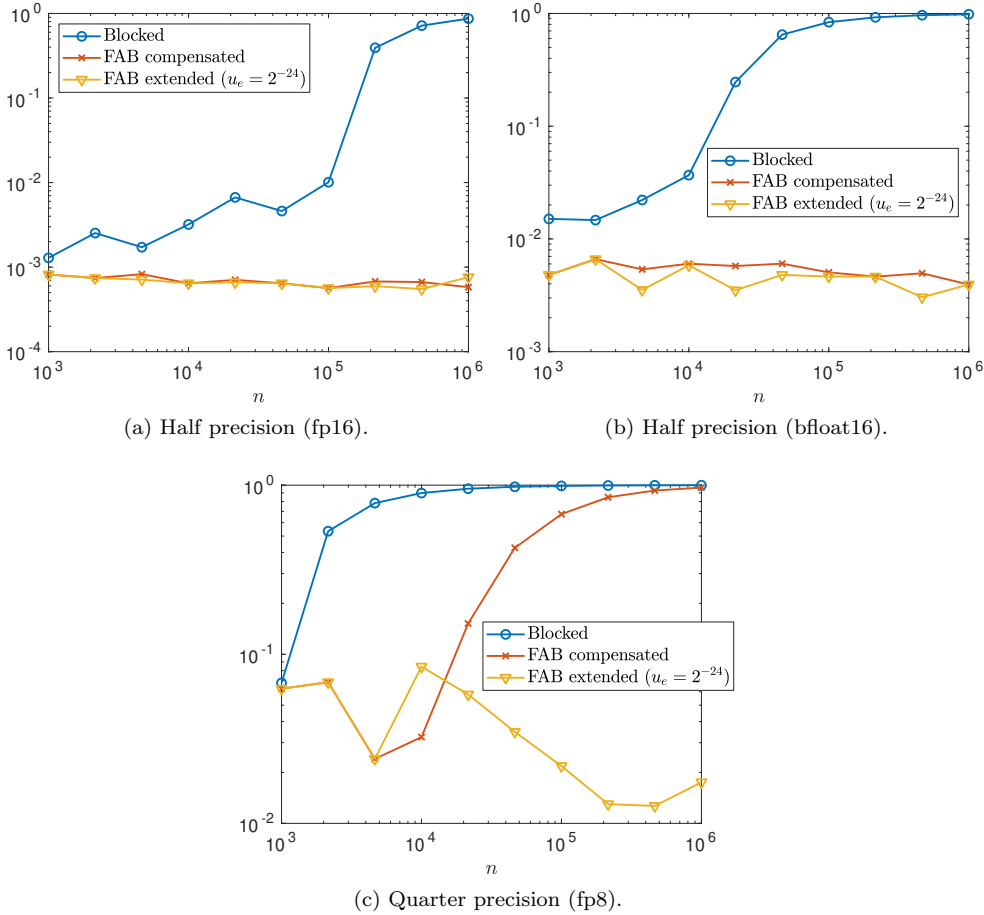


Fig. 5.2: Backward error of the computed sum $\sum_{i=1}^n x_i$ for a vector x with random uniform $[0, 1]$ entries and block size $b = 32$. The methods are blocked summation and `FABsum` with compensated summation or recursive summation in extended precision ($u_e = 2^{-24}$) for `AccurateSum`.

$2b/u_e \approx 1.1 \times 10^9$. For larger n , one could use an even higher extended precision, such as double precision. In fact, provided that stagnation does not occur within local sums, `FABsum` with extended precision summation is able to avoid stagnation as long as standard blocked summation in precision u_e , at a fraction of the cost. For the fp16 and bfloat16 precisions, even though the range of sizes n are too small to show the full effect of stagnation, the variant with extended precision is again noticeably and consistently more accurate than the variant with compensated summation. In contrast to the fp8 case, it is not clear, though, whether this is due to some mild beginning of stagnation or to the effect of the second-order terms analyzed in section 3.1.2.

We conclude with a classical example affected by stagnation: the evaluation of the harmonic series $s = \sum_{i=1}^n 1/i$. In exact arithmetic, s diverges as $n \rightarrow \infty$, whereas in floating-point arithmetic, s is known to “converge” due to stagnation [14], [20]. In Figure 5.3, we compare the value to which s converges for different precisions and

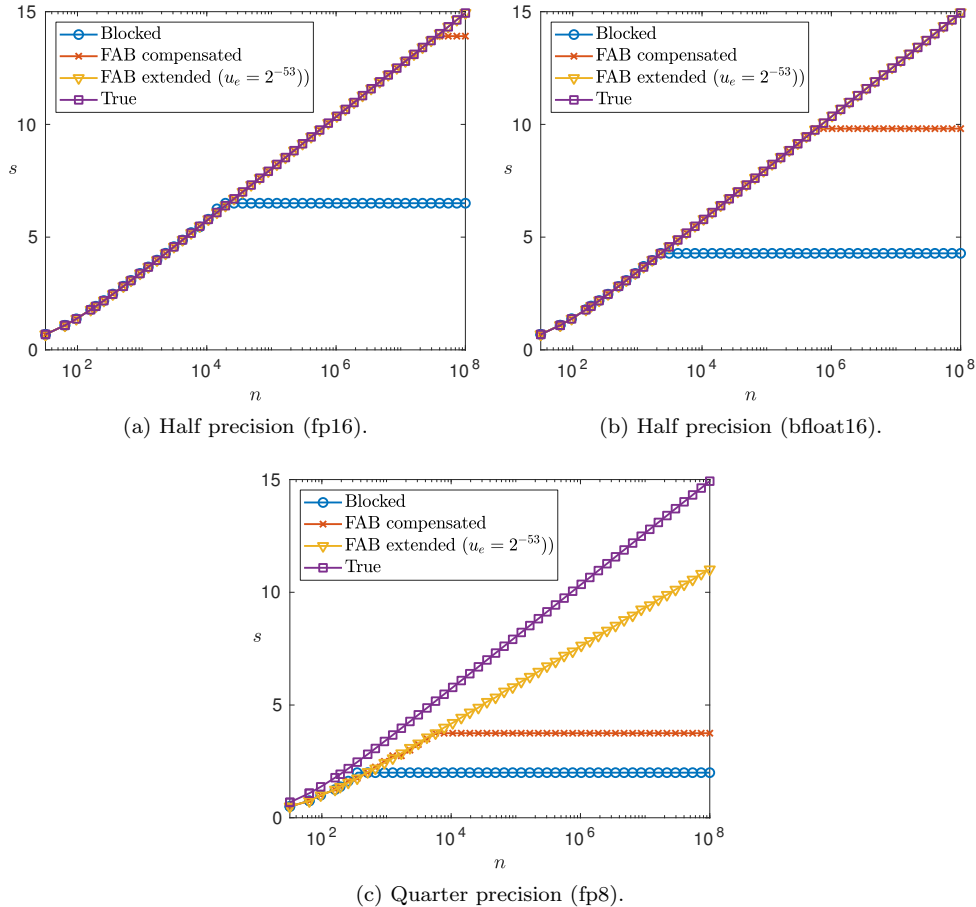


Fig. 5.3: Value of the computed sum $s = \sum_{i=1}^n 1/i$ for different precisions and summation algorithms, using a block size $b = 32$.

summation algorithms. These experiments demonstrate again that `FABsum` with extended precision summation (in this case, double precision) is able to delay stagnation much longer than the other algorithms.

5.2. Matrix multiplication. We now experiment with the state-of-the-art parallel numerical linear algebra library PLASMA [6]. PLASMA includes much of the functionality of the BLAS and LAPACK, but partitions matrices into $t \times t$ blocks called tiles. Independent operations on different tiles are performed concurrently using OpenMP task-based parallelism in order to achieve high performance on shared-memory multicore machines. The tile size is tuned for performance: it must be large enough to allow for computations of high granularity but small enough to expose enough parallelism to feed all the available cores.

By design, PLASMA therefore lends itself naturally to using blocked summation. Note that it is not easy to determine what specific summation algorithm PLASMA uses, because it relies on third-party libraries (in our experiments, Intel MKL) to compute the local operations on tiles. Nevertheless, we expect the default implemen-

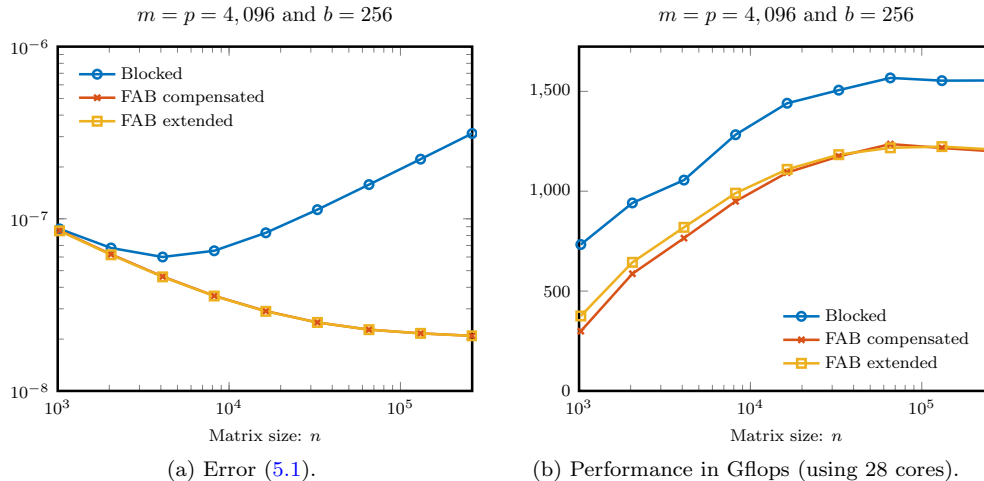


Fig. 5.4: Performance and accuracy of $m \times n$ by $n \times p$ matrix–matrix multiplication in single precision with $m = p = 4096$, block size $b = 256$, tile size $t = 256$, and varying n . We compare the standard implementation based on blocked summation with two versions based on the new **FABsum** algorithm.

tation of the library to suffer from growth of the error with the matrix size. In order to assess how **FABsum** can overcome this issue and at what performance cost, we have implemented **FABsum** in some of the PLASMA routines.²

We experiment with single precision matrix–matrix multiplications $C = AB$, where $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$ are randomly generated matrices with uniform entries in $[0, 1]$. We use two 14-core Intel Broadwell processors on the Kebnekaise supercomputer (Umeå, Sweden). We perform $N_{test} = 10$ consecutive runs and take the maximum performance and the maximum error, measured as

$$(5.1) \quad \frac{\|\widehat{C} - C\|}{\|A\| \|B\|},$$

where the norm is the Frobenius norm and the “exact” C is computed by a double precision multiplication. We wish to study the behavior of the error for large n . Since the matrices would not fit in memory if taken to be square, we fix $m = p = 4096$ and vary n from 10^3 to 2×10^5 .

In Figure 5.4, we compare the performance and accuracy of PLASMA using **FABsum** with the default implementation of the PLASMA library. In this first experiment, we have taken the summation block size b to be equal to the default tile size $t = 256$ because this makes for the most convenient implementation of **FABsum**. Figure 5.4a shows that for the largest matrices, the error is reduced by more than an order of magnitude. In theory, the performance loss should follow the flop overhead, which is about 1% since $b = 256$. However, Figure 5.4b reports a much higher performance loss: while the default implementation achieves a performance peak of over 1500 Gflops for the largest matrices, **FABsum** only reaches about 1200 Gflops, which

²The code is available online at <https://gitlab.com/nla-grp/plasma-17.1-xp>.

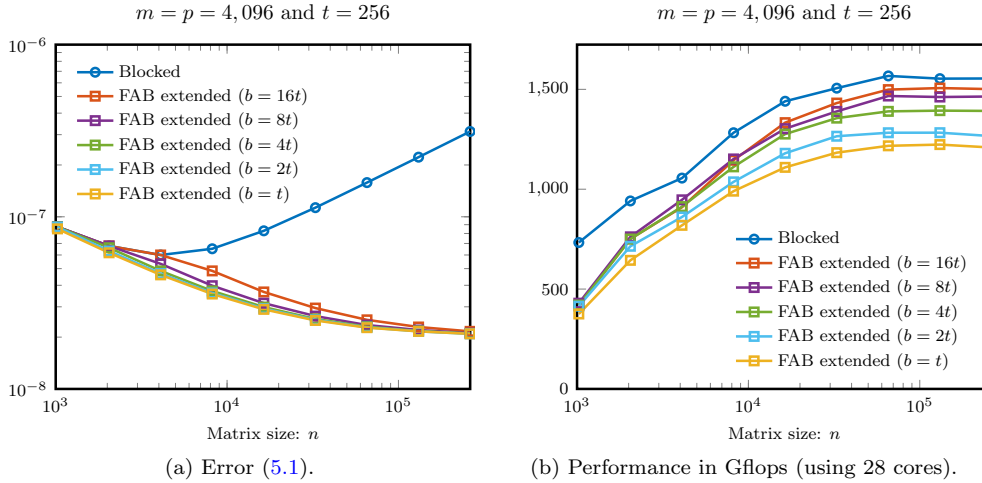


Fig. 5.5: Performance and accuracy of $m \times n$ by $n \times p$ matrix–matrix multiplication in single precision with $m = p = 4096$, tile size $t = 256$, and varying n . We compare the standard implementation based on blocked summation with the FABsum algorithm for different block sizes b varying from t to $16t = 4096$.

represents a 20% loss. This is probably explained by the fact that the extra work required by FABsum cannot achieve a high speed, being of BLAS-2 type rather than BLAS-3.

Therefore, even though it is convenient to take the summation block size b and the tile size t to be the same, doing so does not allow for an entirely satisfying performance/accuracy tradeoff. While dissociating b and t slightly complicates the implementation, it allows us to find the best possible compromise between performance and accuracy, as shown in Figure 5.5. Figure 5.5a shows that increasing b from $t = 256$ to $16t = 4096$ barely impacts the error, even though the bound is 16 times larger (this is partly explained by the probabilistic arguments discussed in the previous section, which suggest the error should increase by only a factor around $\sqrt{16}$ —but even this prediction turns out to be pessimistic for large n). In turn, taking a larger b significantly improves the performance of FABsum: Figure 5.5b shows that with $b = 16t$ the performance loss reduces significantly to about 3%.

6. Conclusion. As problem sizes n in scientific computing continue to grow, rounding error bounds proportional to nu , where u is the unit roundoff, become less satisfactory—especially in the context of the low precision arithmetics that are increasingly available in hardware. Error bounds for (blocked recursive) summation have the form $nu + O(u^2)$ and so can readily exceed 1 in practice. We have proposed a class of blocked summation algorithms, FABsum, that uses a fast algorithm to sum blocks of b numbers and an accurate algorithm to sum the resulting local sums. When the fast algorithm is recursive summation and the accurate algorithm is either compensated summation or recursive summation in extended precision, FABsum has a backward error bound $(b+1)u + O(u^2)$; this bound is about a factor n/b smaller than that for standard summation and, for constant b , does not grow with n to first order.

When FABsum is used in the core linear algebra operations of matrix multiplication

and solution of linear systems by LU factorization, error bounds with constant first order term are again obtained, as shown in section 4. Our numerical experiments show that FABsum does indeed produce substantial reductions to the actual backward errors in practice, and can do so with only a small reduction in performance in the context of the high performance state-of-the-art linear algebra library PLASMA.

FABsum provides an attractive way to obtain more accurate sums and more accurate linear algebra kernel evaluations, especially for low precisions.

Acknowledgments. This research made use of the resources of High Performance Computing Center North (HPC2N). We thank Claude-Pierre Jeannerod and the referees for their helpful suggestions.

REFERENCES

- [1] M. BADIN, L. BIC, M. DILLENCOURT, AND A. NICOLAU, *Improving accuracy for matrix multiplications on GPUs*, Scientific Programming, 19 (2011), pp. 3–11, <https://doi.org/10.3233/SPR-2011-0315>, <https://www.hindawi.com/journals/sp/2011/417569/abs/>.
- [2] E. CARSON AND N. J. HIGHAM, *Accelerating the solution of linear systems by iterative refinement in three precisions*, SIAM J. Sci. Comput., 40 (2018), pp. A817–A847, <https://doi.org/10.1137/17M1140819>.
- [3] A. M. CASTALDO, R. C. WHALEY, AND A. T. CHRONOPOULOS, *Reducing floating point error in dot product using the superbloc family of algorithms*, SIAM J. Sci. Comput., 31 (2008), pp. 1156–1174, <https://doi.org/10.1137/070679946>.
- [4] S. COLLANGE, D. DEFOUR, S. GRAILLAT, AND R. IAKYMCHUK, *Numerical reproducibility for the parallel reduction on multi- and many-core architectures*, Parallel Comput., 49 (2015), pp. 83–97, <https://doi.org/10.1016/j.parco.2015.09.001>.
- [5] J. W. DEMMEL AND H. D. NGUYEN, *Parallel reproducible summation*, IEEE Trans. Automat. Control, 64 (2015), pp. 2060–2070, <https://doi.org/10.1109/TC.2014.2345391>.
- [6] J. DONGARRA, M. GATES, A. HAIDAR, J. KURZAK, P. LUSZCZEK, P. WU, I. YAMAZAKI, A. YARKHAN, M. ABALENKOV, N. BAGHERPOUR, S. HAMMARLING, J. ŠÍSTEK, D. STEVENS, M. ZOUNON, AND S. D. RELTON, *PLASMA: Parallel linear algebra software for multicore using OpenMP*, ACM Trans. Math. Softw., 45 (2019), pp. 16:1–16:35, <https://doi.org/10.1145/3264491>.
- [7] D. GOLDBERG, *What every computer scientist should know about floating-point arithmetic*, ACM Computing Surveys, 23 (1991), pp. 5–48, <https://doi.org/10.1145/103162.103163>.
- [8] A. HAIDAR, S. TOMOV, J. DONGARRA, AND N. J. HIGHAM, *Harnessing GPU tensor cores for fast FP16 arithmetic to speed up mixed-precision iterative refinement solvers*, in Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis, SC '18 (Dallas, TX), Piscataway, NJ, USA, 2018, IEEE Press, pp. 47:1–47:11, <https://doi.org/10.1109/SC.2018.00050>.
- [9] J. R. HAUSER, *Handling floating-point exceptions in numeric programs*, ACM Trans. Program. Lang. Syst., 18 (1996), pp. 139–174, <https://doi.org/10.1145/227699.227701>, <http://doi.acm.org/10.1145/227699.227701>.
- [10] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, second ed., 2002, <https://doi.org/10.1137/1.9780898718027>.
- [11] N. J. HIGHAM AND T. MARY, *A new approach to probabilistic rounding error analysis*, SIAM J. Sci. Comput., 41 (2019), pp. A2815–A2835, <https://doi.org/10.1137/18M1226312>.
- [12] N. J. HIGHAM AND T. MARY, *A new preconditioner that exploits low-rank approximations to factorization error*, SIAM J. Sci. Comput., 41 (2019), pp. A59–A82, <https://doi.org/10.1137/18M1182802>.
- [13] N. J. HIGHAM AND T. MARY, *Sharper probabilistic backward error analysis for basic linear algebra kernels with random data*, MIMS EPrint 2020.4, Manchester Institute for Mathematical Sciences, The University of Manchester, UK, Jan. 2020, <http://eprints.maths.manchester.ac.uk/2743/>.
- [14] N. J. HIGHAM AND S. PRANESH, *Simulating low precision floating-point arithmetic*, SIAM J. Sci. Comput., 41 (2019), pp. C585–C602, <https://doi.org/10.1137/19M1251308>.
- [15] *IEEE Standard for Floating-Point Arithmetic*, IEEE Std 754-2019 (Revision of IEEE 754-2008), IEEE Computer Society, New York, 2019, <https://doi.org/10.1109/IEEESTD.2019>.

- 8766229.
- [16] C.-P. JEANNEROD AND S. M. RUMP, *Improved error bounds for inner products in floating-point arithmetic*, SIAM J. Matrix Anal. Appl., 34 (2013), pp. 338–344, <https://doi.org/10.1137/120894488>.
 - [17] W. KAHAN, *Further remarks on reducing truncation errors*, Comm. ACM, 8 (1965), p. 40, <https://doi.org/10.1145/363707.363723>.
 - [18] D. E. KNUTH, *The Art of Computer Programming, Volume 2, Seminumerical Algorithms*, Addison-Wesley, Reading, MA, USA, third ed., 1998.
 - [19] M. LANGE AND S. M. RUMP, *Sharp estimates for perturbation errors in summations*, Math. Comp., 88 (2019), pp. 349–368, <https://doi.org/10.1090/mcom/3355>.
 - [20] D. MALONE, *To what does the harmonic series converge?*, Irish Math. Soc. Bulletin, 71 (2013), pp. 59–66, <https://www.maths.tcd.ie/pub/ims/bull71/>.
 - [21] C. B. MOLER, *“Half precision” 16-bit floating point arithmetic*. <http://blogs.mathworks.com/cleve/2017/05/08/half-precision-16-bit-floating-point-arithmetic/>, May 2017.
 - [22] O. MÖLLER, *Quasi double-precision in floating point addition*, BIT, 5 (1965), pp. 37–50, <https://doi.org/10.1007/BF01975722>, <http://link.springer.com/10.1007/BF01975722>.
 - [23] W. OETTLI AND W. PRAGER, *Compatibility of approximate solution of linear equations with given error bounds for coefficients and right-hand sides*, Numer. Math., 6 (1964), pp. 405–409, <https://doi.org/10.1007/BF01386090>.
 - [24] S. M. RUMP, *Error estimation of floating-point summation and dot product*, BIT, 52 (2012), pp. 201–220, <https://doi.org/10.1007/s10543-011-0342-4>.