



**HAL**  
open science

## Shadow computation with BFloat16 to estimate the numerical accuracy of summations

David Defour, Pablo de Oliveira Castro, Matei Istioan, Eric Petit

► **To cite this version:**

David Defour, Pablo de Oliveira Castro, Matei Istioan, Eric Petit. Shadow computation with BFloat16 to estimate the numerical accuracy of summations. IEEE 28th Symposium on Computer Arithmetic (ARITH), Jun 2021, Virtual Conference, France. hal-03159965v2

**HAL Id: hal-03159965**

**<https://hal.science/hal-03159965v2>**

Submitted on 26 May 2021

**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.

# Shadow computation with BFloat16 to estimate the numerical accuracy of summations

David Defour  
LAMPS, Univ. of Perpignan  
France.  
Email: david.defour@univ-perp.fr

Pablo de Oliveira Castro  
and Matei Iştoan  
Université Paris-Saclay, UVSQ, LI-PaRAD, ECR  
France  
Email: {pablo.oliveira, matei.istoan}@uvsq.fr

Eric Petit  
Intel Corp.  
Email: eric.petit@intel.com

**Abstract**—In this article, we propose to exploit the new computational capability offered by the Bfloat16 representation format to perform shadow computations and compute estimations of the relative error. We demonstrate and evaluate the assumptions under which shadow computation is valid for the summation problem.

## I. INTRODUCTION

With the finite precision arithmetic available in today’s processors, results are subject to numerical errors which can come from rounding errors. It is difficult to estimate the numerical error of a computation without further investigation, or, when comparing two results, find the one which exhibits the smallest relative error.

Meanwhile, we are witnessing a generalization in the support of low-precision floating-point formats in order to address new use-cases (ex: neural networks). Among them, Bfloat16 (or BF16) leads to a tremendous increase in processing power (8-32x peak FLOPS higher than Binary32). Furthermore, the Intel AMX implementation [1] can be used concurrently with the AVX FP32 units, which means that with careful software design, a shadow BF16 execution can have a low overhead. The BF16 FMA operation uses BF16 inputs and a Binary32 accumulator that can be saved to memory in either Binary32 or BF16. It has been demonstrated that the 32-bit accumulator is a low-overhead, critical feature to handle AI workloads [11], by avoiding absorption (swamping effect)

In this article, we are interested, regarding Binary32 and Bfloat16 computation for the summation problem, in:

- Predicting a result’s accuracy and detect potentially faulty execution using estimators.
- Given two computed results in Binary32, determining the one with the smallest relative error.

## II. BIBLIOGRAPHY AND RELATED WORKS

The study of error and in particular of round-off error due to the use of floating-point arithmetic is an important field of numerical analysis. A *posteriori* error analysis take into account uncertainty and rounding error along with the main computation leading to a computational and memory overhead, slowing the observed computation.

**Funding:** This work was supported by the ANR-20-CE46-0009 InterFLOP

Stochastic arithmetic estimates the error distribution with a stochastic simulation of round-off effects, either applying random rounding at each step of the chain of computation [18] or through Monte-Carlo arithmetic [15]. It can help localizing numerical errors in post-mortem numerical analysis but is impractical for online error estimation.

Interval arithmetic is used to model uncertainty on input data, along with round-off error, in order to guarantee the results [12]. It has a contained memory and computational overhead, but suffers from the dependency problem potentially leading to large uninformative interval.

Error analysis on a sequence of operations can be done by devoting extra memory and computational resources to recompute with extra-accuracy, either by using dedicated arithmetic (ex: MPFR[4], double-double[7]), or compensation [5].

For each of the previously mentioned methods, beside the fact that the overhead in terms of resources can be large, none of them is able to provide numerical information, for an additional cost that is lower than the original computation.

The application range of BF16 has been limited to specific fields, such as AI, due to the limited accuracy of the format. Recently, Henry et al.[6] leveraged BF16 to speed-up computations by splitting Binary32 operations into sub-operations and provide an accuracy analysis. They also propose to use iterative refinement with BF16 as a pre-conditioner to accelerate convergence in a LU solver.

The concept of replicating a sequence of computations, also called shadow replication, has been used in the context of fault tolerance for green cloud computing [14]. The replication can be limited to intermediate data in the context of shadow memory, used for example for numerical debugging [13].

## III. BACKGROUND

Our objective consists of using the extra computational capabilities offered by BF16, available in today’s processors, to estimate the numerical error for the classical (non-compensated) summation problem. In this section, we present the BF16 format and how errors can be bounded.

### A. Bfloat16

Bfloat16 corresponds to a 16-bit shortened version of the IEEE-754 32-bit single-precision floating-point format (i.e. Binary32). The 16 least-significant bits of Binary32 correspond

solely to mantissa bits, which allows for fast conversion to and from Bfloat16. The description of the bit fields of both formats, Bfloat16 and Binary32, are given in Figure 1.

As we can observe, the difference between both formats lies in the fraction field, which is 7-bits long, as opposed to 23-bits in Binary32. This means that the representation range of both formats is almost equivalent.

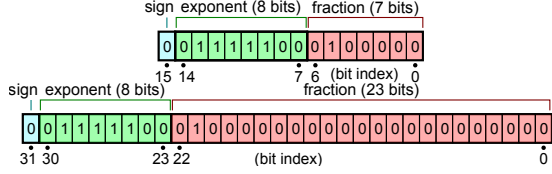


Fig. 1: The Bfloat16 and Binary32 format

### B. Rounding error analysis

We denote by  $\Delta_f(x)$  and  $\bigcirc_f(x)$  the rounding away from 0 and the rounding to nearest of  $x \in \mathbb{R}$ , respectively, in the format  $f$ . Among the considered formats, in this article we use Binary64, Binary32 and Bfloat16, which are represented in the rest of the paper by  $b64$ ,  $b32$ ,  $bf16$ .

In the interval of representable numbers, the output of floating-point operations is impacted by rounding, absorption and catastrophic cancellation errors. In the case of the sum  $s_n = \sum_{i=1}^n x_i$  of a vector  $x_i$  of floating-point numbers, the real absolute error depends on the order of the partial summations. However, there exist metrics which are independent of the order. In rounding to nearest, the approximated computation  $\widehat{s}_n$  of the true sum, can be bounded [8], independently of the order, for all summation methods without overflow as follows:

$$|\widehat{s}_n - s_n| \leq (n-1)u_f \sum_{i=1}^n |x_i| + O(u_f^2) \quad (1)$$

with  $u_f$  the machine precision ( $u_f = 2^{-53}$ ,  $2^{-24}$ ,  $2^{-8}$  for  $b64$ ,  $b32$  and  $bf16$  respectively).

Rump introduced the unit in the first place concept ( $ufp$ ) to perform elegant error analysis, such as in [16]. He proposed algorithms to compute estimations of the error for the summation and dot product problems. For the summation, the algorithm is:

```

 $\widehat{s}_1 = x_1; \widehat{T}_1 = |x_1|$ 
for  $\{k = 2; k \leq n; k++\}$  do
   $\widehat{s}_k = \bigcirc_f(\widehat{s}_{k-1} + x_k)$ 
   $\widehat{T}_k = \bigcirc_f(\widehat{T}_{k-1} + |x_k|)$ 
end for

```

The computed estimator  $\widehat{T}_n$  corresponds to the accumulation of the absolute value of the  $x_k$  for each step of the algorithm. Both  $\widehat{s}_k$  and  $\widehat{T}_k$  must use the same floating-point format. It can be used to derive a computable absolute error bound:

$$\begin{aligned} |\widehat{s}_n - \sum_{i=1}^n x_i| &\leq (n-1) \cdot u_f \cdot ufp(\widehat{T}_n) \\ &\leq (n-1) \cdot u_f \cdot \widehat{T}_n \end{aligned}$$

This bound is known to be a pessimistic upper bound, especially when dealing with random data. Except for very particular cases, the real bound is lower [8]. Using probabilistic assumptions about the rounding error, a relaxed bound can be derived [9] which grows as  $\sqrt{n} \cdot u_f$  instead of  $n \cdot u_f$  for the worst case. In [10], the authors extended their work by proposing a bound depending on the data range distribution independently of  $n$ .

## IV. PROPOSED STRATEGIES

Relative error, or approximation error is a useful metric to estimate the number of correct digits in a result. The question which we address is: how can we provide a computable relative error for a small overhead?

### A. Reference relative error

In order to establish a fair comparison, we have to set a reference value, that we name  $E_{ref}$ . Computing the relative error would require, in the case of the summation problem, computing the reference value using a long accumulator in order to avoid rounding errors, such as the one proposed in [3]. However, for problems of reasonable size and condition number (eg.  $N < 2^{20}$  and  $C < 2^{30}$ ), computing the reference value using Binary64 can be considered enough:

$$E_{ref} = \frac{|\widehat{s}_{b32} - \widehat{s}_{b64}|}{|\widehat{s}_{b64}|} \quad (2)$$

with  $\widehat{s}_f$  corresponding to the computed sum in format  $f$  with rounding to nearest mode.

### B. Mixed bound

Rump's estimator ( $\widehat{T}_n$ ) is only valid when both  $\widehat{T}_n$  and  $\widehat{s}_n$  are computed with identical floating-point formats and rounding modes. It is straightforward to use BF16 to compute  $\widehat{T}_n$ , and only requires few modifications. The proof of Rump's estimator is based on the monotonicity of rounding ([16], equation (3.9)) which states that if  $|x| \leq X$  and  $|y| \leq Y$  then  $|\bigcirc_f(x+y)| \leq \bigcirc_f(X+Y)$  for  $x, y, X, Y \in \mathbb{R}$ .

We propose to alter Rump's algorithm, to replace the computation of  $\widehat{T}_n$  by  $\widehat{B}_n$  computed as follow:

$$B_k = \widehat{B}_{(k-1)} + \Delta_{bf16}(|p_k|); \widehat{B}_k = \bigcirc_{b32}(B_k); \quad (3)$$

This corresponds to loading BF16 approximations of  $|p_k|$ , using rounding away from zero, accumulated in a Binary32 accumulator using rounding to nearest. This corresponds to what is usually available in BF16 hardware.

**Theorem IV.1.** *If  $|x| \leq X$  and  $|y| \leq Y$ , then  $|\bigcirc_{b32}(x+y)| \leq \bigcirc_{b32}(\Delta_{bf16}(X) + \Delta_{bf16}(Y))$  for  $x, y, X, Y \in \mathbb{R}$ .*

*Proof.* Rounding away from zero lets us have  $|x| \leq \Delta_{bf16}(|x|)$ . Upscaling does not involve rounding error, therefore  $\bigcirc_{b32}(\Delta_{bf16}(x)) = \Delta_{bf16}(x)$ . Thanks to the monotonicity of rounding, we get  $|\bigcirc_{b32}(x+y)| \leq \bigcirc_{b32}(\Delta_{bf16}(X) + \Delta_{bf16}(Y))$ .  $\square$

Note that it is possible to consider performing everything in BF16 with rounding away from zero (including the accumulation), based on the fact that double rounding issues do not occur with directed rounding. This leads to  $|\bigcirc_{b32}(x + y)| \leq \bigcirc_{b32}(\Delta_{bf16}(X) + \Delta_{bf16}(Y)) \leq \Delta_{bf16}(\Delta_{bf16}(X) + \Delta_{bf16}(Y))$ . However, an inevitable drift from the real value happens, as it is usually observed with interval arithmetic, but emphasized here by the poor accuracy provided by BF16. We propose the following mixed bound  $E_{mixed}$ , based on Rump's estimator for the absolute value, using BF16 input values:

$$E_{mixed} = (n - 1) \cdot u_{b32} \cdot \frac{\widehat{B_{bf16}}}{|\widehat{s_{b64}}|} \quad (4)$$

with  $\widehat{s_{b64}}$  computed in rounding to nearest and  $\widehat{B_{bf16}}$  computed as in equation 3.

### C. Fully computed bound

Equation 4 cannot be considered a fully computable bound as it requires  $|\widehat{s_{b64}}|$ . We have to bound  $|\widehat{s_{b64}}|^{-1}$  using the computed quantities ( $\widehat{s_{b32}}$  and  $\widehat{B_{bf16}}$ ). One may use the fact that  $|S| \geq |\widehat{s_{b32}}| - (n - 1)u_{b32}\widehat{B_{bf16}}$  to propose the following *computed* estimator:

$$E_{comp} = (n - 1) \cdot u_{b32} \cdot \frac{\widehat{B_{bf16}}}{|\widehat{s_{b32}}| - (n - 1)u_{b32}\widehat{B_{bf16}}} \quad (5)$$

This estimator is valid only if  $|\widehat{s_{b32}}| > (n - 1)u_{b32}\widehat{B_{bf16}}$ .

### D. Approximate estimator number 1

Equation 5 is a valid upper bound corresponding to the worst case scenario. As for other estimators, we can consider a more plausible one. In particular, we use the approximation  $|S| \approx |\widehat{s_{b32}}|$  to propose the  $E_{approx}$  estimator:

$$E_{approx} = (n - 1) \cdot u_{b32} \cdot \frac{\widehat{B_{bf16}}}{|\widehat{s_{b32}}|} \quad (6)$$

## V. RESULTS

This section illustrates through some experiments how BF16 shadow computation can add numerical information on a main computation done in Binary32. To do so, we compare the four estimators on the summation problem. We have executed the naive summation for 5000 vectors of 400 floating-point numbers with condition numbers ranging from  $2^6$  and  $2^{50}$ . We use a random generator inspired by the GenSum function in [17]. Results are reported in Figure 2, which represents the relative error in log scale according to the condition number.

We observe that all four estimators are above (overestimation) the reference error (in blue) for problems of condition number less than  $2^{24}$ . For problems of condition numbers greater than this bound, we observe that the mixed estimator follows through a straight line the real error. The bound  $E_{comp}$  defined in equation 5 (in grey) presents a singularity for problems of condition number greater than  $2^{14}$ . This is due to the validity condition of the estimator ( $|\widehat{s_{b32}}| > (n - 1)u_{b32}\widehat{B_{bf16}}$ ). Finally, the approximate estimator (yellow) behaves as the

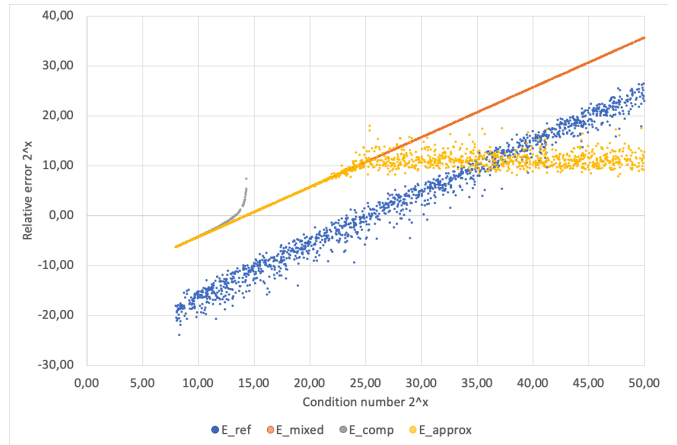


Fig. 2: Relative error for random vectors of size 400 for the four estimators according to the condition number.

mixed estimator for problems of condition numbers less than  $2^{24}$ . With larger condition numbers, it saturates at an error  $2^{10}$ . Even though this estimator is wrong in this case, it still indicates that all bits are lost in the result, which is an important information.

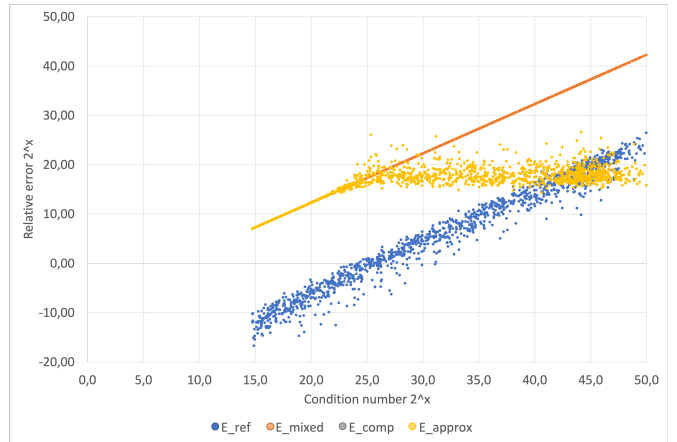


Fig. 3: Relative error for random vectors of size 40000 for the four estimators according to the condition number.

We have conducted tests on vector or larger size (Fig. 3). Previous conclusions stay relevant, we can observe that the differences for each estimators lies in the  $n$  factor (vector size parameter) leading to an upward shift of each plots.

### A. Comparing numerical quality

In this section, we illustrate how the approximated estimator  $E_{approx}$  can be used to compare the numerical quality of two computed results. This could be useful, for example to detect the rise of computer zero during pivoting in LU factorisation [2] and make better selection. We have generated  $10^5$  random vectors of size 400. We counted the number of times the approximated estimator makes a correct prediction to determine which of two results presents the smallest relative error. Results are reported in Figure 4. It represents the

	7	9	11	13	15	17	19	21	23	25	27	29	31	33	35	37	39	41	43	45	47	49	
7	55	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
9	66	58	80	94	98	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
11	100	85	56	83	95	97	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
13	100	97	88	49	81	93	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
15	100	100	97	86	60	84	95	97	98	100	100	100	100	100	100	100	100	100	100	100	100	100	100
17	100	100	98	98	90	59	75	96	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
19	100	100	100	100	95	87	63	77	93	100	100	98	100	100	100	100	100	100	100	100	100	100	100
21	100	100	100	100	98	98	73	56	81	96	100	97	100	100	100	100	100	100	100	100	100	100	100
23	100	100	100	100	98	100	91	79	60	77	77	91	90	89	81	94	88	91	87	93	89	95	95
25	100	100	100	100	100	100	96	98	74	23	50	45	51	45	46	58	62	55	56	58	46	45	45
27	100	100	100	100	100	100	98	100	87	46	19	30	43	43	59	52	54	46	43	39	48	52	52
29	100	100	100	100	100	100	100	100	92	46	31	1	27	35	38	52	51	49	50	55	58	55	55
31	100	100	100	100	100	100	100	100	96	51	48	28	5	30	36	55	58	43	46	54	45	42	42
33	100	100	100	100	100	100	100	100	95	64	43	45	33	3	27	40	40	55	42	46	44	50	50
35	100	100	100	100	100	100	100	100	94	51	54	61	46	36	0	27	37	47	40	55	59	43	43
37	100	100	100	100	100	100	100	100	83	48	58	43	45	46	29	14	20	43	53	52	51	54	54
39	100	100	100	100	100	100	100	100	86	51	42	41	54	44	43	28	3	22	42	41	48	59	59
41	100	100	100	100	100	100	100	100	88	58	44	53	41	48	51	52	28	6	22	43	58	53	53
43	100	100	100	100	100	100	100	100	88	68	45	51	54	78	59	47	51	29	5	28	43	45	45
45	100	100	100	100	100	100	100	100	84	70	36	47	52	46	50	60	46	44	28	8	28	51	51
47	100	100	100	100	100	100	100	100	89	48	44	50	55	42	52	49	49	51	47	8	7	30	30
49	100	100	100	100	100	100	100	100	86	52	55	52	38	56	54	46	52	52	48	45	22	3	3

Fig. 4: Percentage of time that BF16 shadow computation make a correct prediction when comparing the relative error of two summation problems with condition number varying between  $2^6$  and  $2^{50}$ .

percentage of correct predictions according to the condition number of both vectors, indicated in xy-coordinate.

We observe that the estimator correctly predicts which result is the most accurate as long as the difference is significant (diagonal) and the condition number for at least one vector is less than  $2^{23}$  (lower right square). The difference has to be significant between the two vectors for problem of condition number less than  $2^{23}$  as this corresponds to a sharp estimator (yellow line in fig. 2) and a reference relative error that may be better than expected (wide set of blue dots for a given condition number). For the lower right square, this corresponds to the place where the estimator predicts that all bits of accuracy are lost and cannot discriminate.

## VI. CONCLUSION

In this article we proposed a computable estimator for the relative error in the case of the summation problem. This illustrates how BF16 computation resources can be used to add information on numerical quality to a Binary32 computation. Through tests, we have shown that the estimator follows the optimal bound for problems of condition number less than the number of bits available in the mantissa. When the condition number is larger, the estimator remains useful by indicating that all bits of information are lost. Finally, we have shown how the estimator successfully predicts which of two computed results is the most accurate, as long as the difference is significant. As future work, we plan to improve the theoretical ground of the proposed estimator and investigate if it attenuates the rise of computer zero during pivoting in LU factorisation.

## REFERENCES

[1] Intel Architecture Instruction Set Extensions and Future Features Programming Reference - 319433-042 December 2020.  
[2] J. M. Chesneau. The equality relations in scientific computing. *Num. Algo.*, 7:129–143, 1994.

[3] S. Collange, D. Defour, S. Graillat, and R. Iakymchuk. Numerical reproducibility for the parallel reduction on multi- and many-core architectures. *Parallel Computing*, 49:83–97, 2015.  
[4] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier, and P. Zimmermann. Mpfr: A multiple-precision binary floating-point library with correct rounding, 2005.  
[5] S. Graillat, P. Langlois, and N. Louvet. Algorithms for accurate, validated and fast polynomial evaluation. October 01 2009.  
[6] G. Henry, Ping Tak Peter Tang, and A. Heinecke. Leveraging the bfloat16 artificial intelligence datatype for higher-precision computations. In *2019 IEEE 26th Symposium on Computer Arithmetic (ARITH)*, pages 69–76. IEEE, 2019.  
[7] Y. Hida, X. S. Li, and D. H. Bailey. Algorithms for quad-double precision floating point arithmetic. In *Proc. 15th IEEE Symposium on Computer Arithmetic*, pages 155–162. IEEE Computer Society Press, Los Alamitos, CA, USA, 2001.  
[8] N. Higham. The accuracy of floating point summation. *SIAM J. Sci. Comput.*, 14:783–799, 1993.  
[9] N. J. Higham and T. Mary. A new approach to probabilistic rounding error analysis. *SIAM J. Sci. Comput.*, 41(5):A2815–A2835, 2019.  
[10] N. J. Higham and T. Mary. Sharper probabilistic backward error analysis for basic linear algebra kernels with random data. *SIAM Journal on Scientific Computing*, 42(5):A3427–A3446, 2020.  
[11] D. Kalamkar et al. A Study of BFLOAT16 for Deep Learning Training, 2019.  
[12] U. W. Kulisch. Mathematics and speed for interval arithmetic: A complement to iee 1788. *ACM Trans. Math. Softw.*, 45(1):5:1–5:22, 2019.  
[13] M. O. Lam and B. L. Rountree. Floating-point shadow value analysis. In *ESPT@SC*, pages 18–25. IEEE, 2016.  
[14] B. N. Mills, T. Znati, and R. G. Melhem. Shadow computing: An energy-aware fault tolerant computing model. In *ICNC*, pages 73–77. IEEE, 2014.  
[15] D. Stott Parker. Monte Carlo arithmetic: exploiting randomness in floating-point arithmetic. Technical Report CSD 970002, Department of Computer Science, University of California, Los Angeles, Los Angeles, CA, USA, 1997.  
[16] Siegfried M. Rump. Error estimation of floating-point summation and dot product. *BIT Numerical Mathematics*, 52(1):201–220, March 2012.  
[17] S.M. Rump. INTLAB - INTerval LABoratory. In Tibor Csendes, editor, *Developments in Reliable Computing*, pages 77–104. Kluwer Academic Publishers, Dordrecht, 1999. <http://www.ti3.tuhh.de/rump/>.  
[18] J. Vignes. A stochastic arithmetic for reliable scientific computation. *Math. Comp. Simul.*, 35:233–261, 1993.