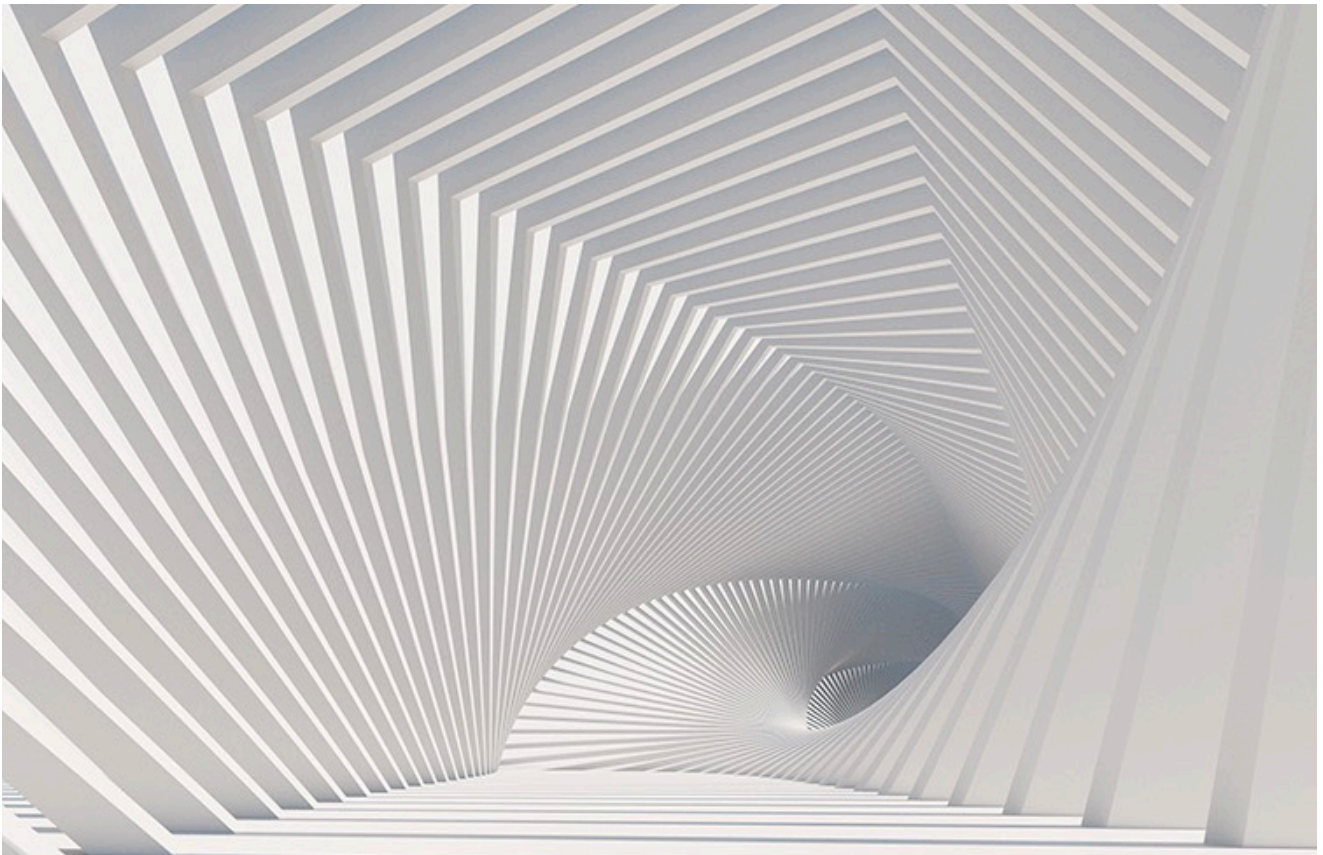


# FLOATING-POINT ERRORS AND COMPENSATION ALGORITHM



Ao XIE

Ingénierie des modèles et Simulation

Institut Supérieur d'Informatique, de  
Modélisation et de leurs Applications

# I ntroduction

As the foundation of modern computing, floating point numbers play an important role in scientific computing, financial modeling and engineering. However, due to the limitations of binary representation, floating-point operations inherently introduce errors and may accumulate over multiple operations.

The representation of floating-point numbers follows the IEEE 754 standard[1], which defines a finite-precision binary floating-point format. For integers, an exact representation is possible based on this definition. However, for decimals, exact representation is not possible. For example, the binary representation of 0.1 is an infinite loop decimal. This limitation requires truncation of decimals, which introduces rounding errors[2].

In addition, floating-point endings have a limited number of bits, which limits the precision of floating-point numbers. When computing values with large magnitude differences, loss of valid digits may occur, known as catastrophic elimination[1].

In this context, this experiment will validate the precision error of floating point numbers themselves and use Kahan Algorithm[3] and its improved algorithm, the Neumaier [4], to try to reduce this error.

## E xperimental Data Set

In total, three different forms of datasets were set up in this experiment to try to reproduce the problems that arise during the computation of floating point numbers.

### 1 Random Dataset

The use of random numbers between 0 and 1 is used to simulate regular numerical distributions. This data is used to simulate scenarios of random number distributions that occur widely in real-world applications.

### 2 Uniform Dataset

Using the same decimals is used to validate that the algorithm handles the accumulation of uniformly distributed data. This data is used to simulate the errors generated during simple accumulation and also for the cases in the course.

### 3 Extreme Dataset

The use of alternating extremely large and extremely small values was used to test the stability of the algorithm when dealing with large spans of values.

# A lgorithms

A total of three different algorithms were used in this experiment to calculate the different errors they produce. These are the use of the Native method as the baseline model and the other two compensation algorithms.

## 1 Native Summation Algorithm

In this experiment, the most primitive cumulative approach was used as the baseline model for the experiment. The role of the baseline model is to verify the errors that can occur in floating-point numbers when computed without correction and whether the correction algorithm can correct these errors to make a reference.

## 2 Kahan Summation Algorithm

Kahan summation algorithm is a computational method proposed by William Kahan in 1965[3] for improving the accuracy of numerical calculations.

The core idea of this algorithm is to record the errors rounded off in each addition operation by introducing a compensation variable, and gradually compensate these errors back in subsequent calculations. The algorithm is shown in Algorithm 1.

---

**Algorithm 1: Kahan Summation**

---

**Input:** Array  $arr$  of size  $n$   
**Output:** Sum of the array with error compensation

```
1 Initialize  $sum \leftarrow 0.0$ ,  $c \leftarrow 0.0$  // Compensation for lost low-order bits
2 for  $i \leftarrow 1$  to  $n$  do
3    $y \leftarrow arr[i] - c$  // Correct the current value with compensation
4    $t \leftarrow sum + y$ 
5    $c \leftarrow (t - sum) - y$  // Update the compensation value
6    $sum \leftarrow t$ 
7 end
8 return  $sum$ 
```

---

This algorithm is effective in compensating rounding errors without introducing a lot of extra operations and has a simple structure for easy completion.

### 3 Neumaier Summation Algorithm

An improved version of the Kahan summation algorithm was proposed by L. Neumaier in 1974[4] to address the stability and robustness of the Kahan summation algorithm in the face of summing over large spans of data.

The core idea of the Neumaier summation algorithm is to complete the compensation in the Kahan summation algorithm by means of dynamic compensation. Specifically, its by determining the relationship between the absolute size of the current sum and the cumulative value, and adjusted to the appropriate way of error compensation. Its specific implementation is shown in Algorithm 2.

---

**Algorithm 2: Neumaier Summation**

---

**Input:** Array *arr* of size *n*

**Output:** Sum of the array with enhanced error compensation

```
1 Initialize sum  $\leftarrow$  0.0, c  $\leftarrow$  0.0           // Enhanced compensation for
   rounding errors
2 for i  $\leftarrow$  1 to n do
3   y  $\leftarrow$  arr[i] - c
4   t  $\leftarrow$  sum + y
5   if  $|sum| \geq |y|$  then
6     c  $\leftarrow$  c + (sum - t) + y
7   else
8     c  $\leftarrow$  c + (y - t) + sum
9   end
10  sum  $\leftarrow$  t
11 end
12 return sum + c
```

---

This algorithm presents stronger robustness error compensation with a slight increase in computational complexity.

## Experimental Record

The main content of this section is to convert the above algorithms into C code and present the final results of the experiments in the form of tables. The experimental results will also be analyzed and summarized to compare the performance, accuracy and applicability of different algorithms in an intuitive way, so as to provide support and reference for the application of the algorithms in practical scenarios.

In addition to the baseline model, the experiments introduced the MPFR library[5] as a high-precision reference standard. In this case, the library is utilized to obtain theoretical results for the ideal case, which can be used as a reference standard for practical calculations.

## 1 Random Dataset

During the computation of random data, each result is computed using 1000 randomly generated random numbers between 0 and 1 and accumulated. The final results are shown in Table 1.

Table 1 - Cumulative test for random data

	MPFR	Native Result	Native Error	Kahan Result	Kahan Error	Neumaier Result	Neumaier Error
1	498.544783810407295732	498.544783810407579949	0.00000000000000284217	498.544783810407295732	0.000000000000000000	498.544783810496539900	0.00000000000089244168
2	479.816517785106100291	479.816517785105872917	0.00000000000000227374	479.816517785106100291	0.000000000000000000	479.816517784974223559	0.00000000000131876732
3	513.348933232179319930	513.348933232178978869	0.00000000000000341061	513.348933232179319930	0.000000000000000000	513.348933232182730535	0.00000000000003410605
4	505.267555140549120551	505.267555140549120551	0.000000000000000000	505.267555140549120551	0.000000000000000000	505.267555140800936897	0.00000000000251816346
5	498.147556402789234653	498.147556402789064123	0.00000000000000170530	498.147556402789234653	0.000000000000000000	498.147556402619215987	0.00000000000170018666
6	506.182040249547924304	506.182040249547753774	0.00000000000000170530	506.182040249547924304	0.000000000000000000	506.182040249561623568	0.00000000000013699264
7	487.139282804047354603	487.139282804047297759	0.000000000000000056843	487.139282804047354603	0.000000000000000000	487.139282804100275825	0.00000000000052921223
8	494.903422127898522831	494.903422127898238614	0.00000000000000284217	494.903422127898522831	0.000000000000000000	494.903422127788530815	0.00000000000109992015
9	475.860664682397953129	475.860664682398294190	0.00000000000000341061	475.860664682397953129	0.000000000000000000	475.860664682248511781	0.00000000000149441348
10	485.702045298508380711	485.702045298508608084	0.00000000000000227374	485.702045298508380711	0.000000000000000000	485.702045298365760573	0.00000000000142620138

In fact, during the experiments in this section we try to use the MT random number generator[6] to get better random numbers for our experiments. But the generator was not used in the end due to the fact that the high resolution MT random numbers are close to uniformly distributed in the interval, a property that largely reduces the accumulation of rounding errors.

Based on the information in the table, it is easy to see that for random numbers between zero and one, the computational error of the baseline model is less than that of Neumaier algorithm when MPFR is used as the reference standard, but it shows less stability and accuracy relative to Kahan algorithm.

## 2 Uniform Dataset

For the uniform data test, since the results were the same for each experiment, data from 0.1 to 1.0 directly cumulative in order of 0.1 were used as the test data. The final test results are shown in Table 2.

Table 2 - Cumulative test for uniform data

	Data	MPFR	Native Result	Native Error	Kahan Result	Kahan Error	Meumaier Result	Meumaier Error
1	0.1	100.0000 0000000 0000000	99.99999 99999985 93125	0.00000 0000001 406875	100.0000 0000000 0000000	0.00000 0000000 0000000	100.0000 0000000 6153300	0.00000 0000006 153300
2	0.2	200.0000 0000000 0000000	199.99999 99999971 86251	0.00000 0000002 813749	200.0000 0000000 0000000	0.00000 0000000 0000000	200.0000 0000001 2306600	0.00000 0000012 306600
3	0.3	300.0000 0000000 0000000	300.0000 0000000 5627498	0.00000 0000005 627498	300.0000 0000000 0000000	0.00000 0000000 0000000	300.0000 0000000 0170530	0.00000 0000000 170530
4	0.4	400.000 0000000 0000000 0	399.9999 9999999 4372502	0.00000 0000005 627498	400.0000 0000000 0000000	0.00000 0000000 0000000	400.0000 0000002 4613200	0.00000 0000024 613200
5	0.5	500.0000 0000000 0000000	500.0000 0000000 0000000	0.00000 0000000 0000000	500.0000 0000000 0000000	0.00000 0000000 0000000	500.0000 0000000 0000000	0.00000 0000000 0000000
6	0.6	600.000 0000000 0000000 0	600.000 0000000 11254997	0.00000 00000112 54997	600.0000 0000000 0000000	0.00000 0000000 0000000	600.0000 0000000 0341061	0.00000 0000000 341061
7	0.7	700.0000 0000000 0000000	700.0000 0000000 6366463	0.00000 0000006 366463	700.0000 0000000 0000000	0.00000 0000000 0000000	700.0000 0000003 2059688	0.00000 0000032 059688

	Data	MPFR	Native Result	Native Error	Kahan Result	Kahan Error	Meumaier Result	Meumaier Error
8	0.8	800.0000 0000000 0000000	799.9999 99999988 745003	0.00000 00000112 54997	800.0000 0000000 0000000	0.00000 0000000 0000000	800.0000 0000004 9226401	0.00000 0000049 226401
9	0.9	900.000 0000000 0000000 0	899.9999 99999984 879651	0.00000 00000151 20349	900.0000 0000000 0000000	0.00000 0000000 0000000	900.0000 0000002 8762770	0.00000 0000028 762770
10	1.0	1000.000 0000000 0000000 0	1000.000 0000000 0000000 0	0.00000 0000000 0000000	1000.000 0000000 0000000 0	0.00000 0000000 0000000	1000.000 0000000 0000000 0	0.00000 0000000 0000000 0000000

According to the data in the table, it is still found that Kahan algorithm performs better for the baseline and Nermaie algorithm. And for decimal numbers that are easy to represent in binary (e.g., 0.5 and 1.0), no rounding error occurs.

### 3 Extreme Dataset

Based on the data obtained earlier, it can be seen that the results obtained by the Neumaier algorithm are not satisfactory, and this result does not have a significant impact in the above conclusions due to the fact that, according to its proposed conception, this algorithm is proposed in order to improve the accuracy when the data variance is large. Therefore in this subsection, an attempt will be made to correct this algorithm. Namely, the algorithm will include a comparison of the size for 1e6 and then two types of compensation will be added. For the experimental results after correcting this algorithm are shown in Table 3.

Table 3 - Cumulative test for extreme data

	MPFR	Native Result	Native Error	Kahan Result	Kahan Error	Neumaier Result	Neumaier Error
1	2944773162 291.7241210 937500000 00	2944773162 291.7426757 812500000 00	0.01855468 75000000 00	2944773162 291.7241210 937500000 00	0.0000000 00000000 000	2944773162 291.6704101 562500000 00	0.05371093 75000000 00
2	4699571328 440.020507 812500000 000	4699571328 439.992187 50000000 0000	0.02832031 25000000 00	4699571328 440.020507 812500000 000	0.0000000 00000000 000	4699571328 440.223632 812500000 000	0.20312500 00000000 00



	MPFR	Native Result	Native Error	Kahan Result	Kahan Error	Neumaier Result	Neumaier Error
<b>3</b>	420804428 8774.87500 00000000 00000	420804428 8774.89257 812500000 0000	0.017578125 00000000 0	420804428 8774.874511 718750000 000	0.0004882 812500000 00	420804428 8775.00830 078125000 0000	0.13330078 125000000 0
<b>4</b>	266736454 5043.21484 37500000 00000	266736454 5043.21289 06250000 00000	0.00195312 50000000 00	266736454 5043.21484 375000000 0000	0.00000000 00000000 000	266736454 5043.182617 187500000 000	0.03222656 25000000 00
<b>5</b>	4274142913 011.098144 531250000 000	4274142913 011.0488281 25000000 000	0.04931640 62500000 00	4274142913 011.0981445 312500000 00	0.00000000 00000000 000	4274142913 011.1655273 437500000 00	0.06738281 25000000 00
<b>6</b>	248769964 9446.86425 781250000 0000	248769964 9446.82373 046875000 0000	0.04052734 375000000 0	248769964 9446.86425 781250000 0000	0.00000000 00000000 000	248769964 9446.99267 578125000 0000	0.128417968 75000000 0
<b>7</b>	4094478017 414.7475585 937500000 00	4094478017 414.668457 031250000 000	0.07910156 25000000 00	4094478017 414.7475585 937500000 00	0.00000000 00000000 000	4094478017 414.663085 937500000 000	0.08447265 62500000 00
<b>8</b>	2717640620 071.469238 281250000 000	2717640620 071.4619140 62500000 000	0.00732421 875000000 0	2717640620 071.469238 281250000 000	0.00000000 00000000 000	2717640620 071.5678710 937500000 00	0.09863281 25000000 00
<b>9</b>	3735716926 05.8669433 593750000 00	3735716926 05.8659667 968750000 00	0.0009765 62500000 000	3735716926 05.8669433 593750000 00	0.00000000 00000000 000	3735716926 05.8804931 640625000 00	0.01354980 468750000 0
<b>10</b>	496396582 5884.91210 937500000 0000	496396582 5884.93164 06250000 00000	0.01953125 00000000 00	496396582 5884.911132 812500000 000	0.0009765 62500000 000	496396582 5885.09179 687500000 0000	0.17968750 00000000 00

According to the experimental results, despite the attempts to change the Neumaier algorithm, the performance results were still unsatisfactory and even performed worse than the baseline model. However, in this experiment, Kahan algorithm still gave relatively good results, and despite some small observable errors in some cases, the overall results were still better than the baseline model and Neumaier algorithm.



# C onclusion

The purpose of this experiment is to compare the results obtained by the native summation method, Kahan summation algorithm and Neumaier summation algorithm.

And to use the results obtained from MPFR library as a reference ideal result for the data.

Based on the experimental results, it can be clearly concluded that the traditional method (i.e., the baseline model) produces a certain rounding error for the vast majority of fractional arithmetic operations from 0 to 1, and therefore this method has absolute limitations. At the same time, it can be seen that Kahan algorithm gets better performance in all three cases listed in the experiment. A small amount of error occurs only when there are large fluctuations in the size of the data.

The result of Neumaier algorithm in the experiment is not a normal result, and according to its proposed purpose, a third experiment was designed to verify its effect. However, for all three experiments, this algorithm did not give an excellent result in terms of error, and even in some experimental cases the error obtained was larger than the baseline model. This phenomenon occurs both before and after the algorithm correction, i.e., the possible problem is in the implementation of the algorithm.

Kahan algorithm continues to have a profound impact in currently popular fields, such as Amro Eldebiky et al. who proposed a technique for error compensation in simulation computation using a similar idea to reduce the impact of hardware noise on neural network training and inference[7], and Yifei Cheng et al. who proposed a technique to reduce communication overheads using a similar approach[8]. Therefore research based on this algorithm as well as its ideas is still of research value at the present time.

# Bibliography

- [1] "754-2019 - IEEE Standard for Floating-Point Arithmetic | IEEE Standard | IEEE Xplore," doi: 10.1109/IEEESTD.2019.8766229.
- [2] GoldbergDavid, "What every computer scientist should know about floating-point arithmetic," *ACM Computing Surveys (CSUR)*, vol. 23, no. 1, 1991-03-01, doi: 10.1145/103162.103163.
- [3] KahanW., "Pracniques: further remarks on reducing truncation errors," *Communications of the ACM*, vol. 8, no. 1, 1965-01-01, doi: 10.1145/363707.363723.
- [4] A. Neumaier, "Rundungsfehleranalyse einiger Verfahren zur Summation endlicher Summen," *Zamm-zeitschrift Fur Angewandte Mathematik Und Mechanik*, 1974, doi: 10.1002/ZAMM.19740540106">10.1002/ZAMM.19740540106</a>.
- [5] FousseLaurent, HanrotGuillaume, LefèvreVincent, PélissierPatrick, and ZimmermannPaul, "MPFR," *ACM Transactions on Mathematical Software (TOMS)*, vol. 33, no. 2, 2007-06-01, doi: 10.1145/1236463.1236468.
- [6] MatsumotoMakoto and NishimuraTakuji, "Mersenne twister," *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, vol. 8, no. 1, 1998-01-01, doi: 10.1145/272991.272995.
- [7] A. Eldebiky, G. L. Zhang, G. Boecherer, B. Li, and U. Schlichtmann, "CorrectNet: Robustness Enhancement of Analog In-Memory Computing for Neural Networks by Error Suppression and Compensation," 2022/11/27, doi: 10.48550/arXiv.2211.14917.
- [8] Y. Cheng et al., "Communication-Efficient Distributed Learning with Local Immediate Error Compensation," 2024/02/19, doi: 10.48550/arXiv.2402.11857.