

# Just-in-Accuracy: Mobile Approach to Uncertainty

Martine Ceberio, Christoph Lauter, and Vladik Kreinovich  
Department of Computer Science  
University of Texas at El Paso  
500 W. University  
El Paso, TX 79968, USA  
mceberio@utep.edu, cqlauter@utep.edu, vladik@utep.edu

## Abstract

To make a mobile device last longer, we need to limit computations to a bare minimum. One way to do that, in complex control and decision making problems, is to limit precision with which we do computations, i.e., limit the number of bits in the numbers' representation. A problem is that often, we do not know with what precision should we do computations to get the desired accuracy of the result. What we propose is to first do computations with very low precision, then, based on these computations, estimate what precision is needed to achieve the given accuracy, and then perform computations with this precision.

## 1 Formulation of the Problem

**For mobile devices, an important restriction is energy.** Mobile devices are very convenient, but they need to be recharged ever so often – e.g., after a certain number of hours. This need comes from the fact a mobile device can store only a limited amount of energy. Every time we perform a computation, we use some energy.

**How to make mobile devices last longer.** Since every bit operation requires energy, the only way to make a mobile device last longer is to reduce the number of bit operations needed to perform the corresponding computations.

The number  $B$  of bit operations can be estimated as the number  $A$  of arithmetic operations times the number  $b$  of bit operations needed for each arithmetic operation:  $B = A \cdot b$ . How can we minimize this product?

The number  $A$  of arithmetic operations is determined by the algorithm we use. Researchers are constantly working on decreasing this amount, so we can safely assume that this number is, at present, as small as we can make it. Thus, the only way to decrease the number of bit operations is to decrease the number of bit operation needed to perform one arithmetic operation.

**How can we decrease the number of bit operations.** For each arithmetic operation, the number of bit operations depends on the number  $n$  of bits used to represent a number. The more bits we need to process to perform each operation, the more bit operations we need. So, the only way to make a mobile device last longer seems to be using fewer bits to represent the corresponding numbers.

**But can we do it?** There is a reason why modern computers use a large number of bits (usually, 64, sometimes 32) to perform arithmetic operations: many computations require high precision, and computations with lower precision result in lower accuracy than we want. For some computations, even 64 bits are not enough, we need double precision (i.e., 128 bits) or even higher.

Because of such computations, we cannot simply reduce the number of bits used to represent each number – that will make many computations impossible.

**So what can we do?** So how can we decrease the overall amount of bit operations and still be able to perform computations requiring high accuracy?

For some algorithms, we know what precision to use to achieve the desired accuracy. For example, in most applications of deep learning, 8-bit computations are sufficient; see, e.g., [4]. In such cases, this is exactly the precision that we need to use for such computations – provided, of course, that the CPU allows computations with different precision.

This knowledge is available for many existing algorithms. However, new algorithm appear all the times, algorithms for which such a study of needed precision has not yet been done. What can we do in this case – other than compute all these algorithm with the highest precision?

**What we do in this paper.** In this paper, we provide a possible solution to this problem: namely, we show how to decrease the overall number of bit operations without sacrificing the desired accuracy.

## 2 What We Propose

**Our main idea.** If we have an algorithm for which we do not know what precision we need to achieve the desired accuracy, then:

- first, we perform computations with low precision;
- based on results of these computations, we determine what precision is absolutely needed to achieve the desired accuracy; and then
- we perform computations with thus determined precision.

*Terminological comment.* In analogy with just-in-time delivery, when we save on storage expenses by scheduling delivery for exactly the time at which the delivered objects are needed, we call the proposed approach – in which we

exactly as many digits as needed to achieve the desired accuracy, no more, no less – *just-in-accuracy* approach.

**How we can implement this idea.** In order to implement this idea, we need to do the following three tasks:

- first, we need to find out how to estimate the accuracy of the result of computations with low precision;
- second, we need to find out how to use this estimate to determine the desired precision;
- finally, we need to make sure that we indeed decrease overall number of bit operations.

Let us consider these tasks one by one.

*Comment.* The first two tasks are clearly needed, but why is the third task important? Because:

- on the one hand, when low-precision computations are already sufficient, the proposed approach definitely decreases the number of bit operations;
- on the other hand, for problems that really need high-precision computations, we have to perform these computations anyway; so in our approach, in addition to these high-precision computations, we also perform additional low-precision computation – and thus, increase the number of bit operations.

We need to make sure that decrease is larger than the increase – then the overall number of bit operations will decrease.

**How to estimate the accuracy of the result of low-precision computations.** The relative inaccuracy caused by limited precision is relatively small: if we use  $n$ -bit precision, then the relative round-off error is of order  $2^{-n}$ . Even if we use a very low 8-bit precision, this error is about  $2^{-8} \approx 1.5\%$ .

Thus, we can use the idea typically used in physics: we expand the dependence of the result of the round-off errors in Taylor series and keep only the first few terms in this expansion; see, e.g., [2, 5]. For the relative error of 1.5%, its square is about 0.02% – which is much smaller than the error itself. Thus, to estimate the effect of these errors, we can safely ignore terms which are quadratic (and of higher order) in terms of these errors and only keep linear terms. So, the accuracy of the resulting computations is a linear function of these errors – i.e., it is proportional to  $2^{-n}$ .

How can we estimate this error? We can repeat the low-precision computations with two different low precisions:  $n_1$  and  $n_2 > n_1$ ; for example, we can take  $n_2 = n_1 + 1$ . For  $n_2 = n_1 + 1$ , the result  $r_1$  of the  $n_1$ -precision computations differs from the unknown actual value  $a$  by some value  $r_1 - a \approx c \cdot 2^{-n_1}$ , while the result  $r_2$  is the  $n_2$ -precision computations differs by the value

$$r_2 - a \approx c \cdot 2^{-(n_1+1)} = 0.5 \cdot c \cdot 2^{-n_1}.$$

Thus, we have

$$r_1 - r_2 = (r_1 - a) - (r_2 - a) \approx c \cdot 2^{-n_1} - 0.5 \cdot c \cdot 2^{-n_1} = 0.5 \cdot c \cdot 2^{-n_1} \approx r_2 - a.$$

Hence, the difference  $r_1 - r_2$  between the results of these two computations can be used as an approximate estimate for the accuracy  $r_2 - a$  of the somewhat-more-precise computation result  $r_2$ . The corresponding relative accuracy is thus approximately equal to the ratio

$$\frac{|r_1 - r_2|}{|r_2|}.$$

It is reasonable to gauge this relative accuracy by the number of correct digits in the binary expansion of the computation result  $r_2$ . Having  $d$  digits means relative accuracy  $2^{-d}$ . Thus, this value  $d$  can be determined from the approximate equality

$$2^{-d} \approx \frac{|r_1 - r_2|}{|r_2|},$$

hence

$$d \approx \log_2(|r_2|) - \log_2(|r_1 - r_2|).$$

**How to determine the desired precision.** Let us denote the desired number of correct bits in the computation result by  $k$ . This means that we want to have the relative accuracy  $2^{-k}$  of the computation result. How many bits do we need to use in our computations to reach this accuracy?

Due to our linearity assumption, in general, the relative accuracy resulting from computations with  $n$  digits – for which the relative round-off error is about  $2^{-n}$  – is approximately equal to  $c \cdot 2^{-n}$  for some constant  $c$ . To find the value of this constant  $c$ , we need to take into account that, based on our low-precision computations, we know that the relative round-off error  $2^{-n_2}$  leads to the accuracy  $2^{-d}$  in the computation result. Thus,  $2^{-d} = c \cdot 2^{-n_2}$ , so  $c = 2^{n_2-d}$ .

We want to find out the value  $n$  for which the resulting accuracy is  $c \cdot 2^{-n} = 2^{-k}$ . Substituting the above formula for  $c$  into this expression, we conclude that  $2^{n_2-d} \cdot 2^{-n} = 2^{n_2-d-n} = 2^{-k}$ . Taking binary logarithm of both sides of this equality, we get  $n_2 - d - n = -k$ , so  $n = k + n_2 - d$ .

Let us summarize what we have found.

**Resulting algorithm.** Suppose that we want to get the result with  $k$  significant digits. Let us select some small number  $n_1$ .

- First, we perform the computations with  $n_1$  bits and with  $n_2 = n_1 + 1$  bits, and get approximate results  $r_1$  and  $r_2$ .
- Then, we find  $d \approx \log_2(|r_2|) - \log_2(|r_1 - r_2|)$ .
- Finally, we perform computations again, this time with  $n = k + n_2 - d$  bits.

### 3 The Proposed Algorithm Does Decrease the Overall Number of Bit Operations

**How can we decide whether this algorithm decreases the overall number of bit operations?** To answer this question, we need to compare the overall number of bit operations in two settings:

- the traditional setting, when all the computations are performed with the same high precision  $N$ , and
- the proposed setting, when we first perform two computations with low-precision values  $n_1$  and  $n_2$ , and then perform computations with the needed precision  $n$ .

To compare these two settings, we need to know:

- how frequent are situations with different values  $n$  of needed precision, and
- how the number of bit operations grows with  $n$ .

Let us try to answer these two questions.

**How frequent are situations with different values  $n$  of needed precision?** As we have mentioned, for different problems, we need different precision, i.e., different numbers  $n$  of bits-per-number, from  $n = 1$  to  $n = N$  for some large  $N$  (e.g.,  $N = 64$  or  $N = 128$ ). We do not know the frequency with which different values  $n$  appear, and we have no reason to believe that some of these values are more frequent than other. Thus, it is reasonable to use Laplace Indeterminacy Principle (see, e.g., [3]) and conclude that all these  $N$  values are equally probable, i.e., that each of them occurs with the same probability  $1/N$ .

**How does the number of bit operations needed for each arithmetic operation change with  $n$ ?**

- For addition and subtraction, this number is proportional to  $n$ : we need a constant number of operations per bit.
- For multiplication, we need order of  $n^2$  operations: indeed, the usual multiplication algorithm means that for each digit in the second number, we use  $n$  bit operations to multiply the first number by this digit, and then we add  $n$   $n$ -digit results.

Since  $n^2 \gg n$ , in the first approximation, we can simply take into account  $n^2$ -operations – e.g., multiplication – and ignore time needed for addition and subtraction. Thus, we can conclude that the number of bit operations grows with  $n$  as  $c \cdot n^2$  for some constant  $c$ .

**Now, we are ready to compare the two settings.** Since we answered both questions, we can now provide the desired comparison.

- In the traditional setting, for each problem, we need  $c \cdot N^2$  bit operations.
- In the proposed setting, for needed precision  $n$ , we need  $c \cdot n_1^2 + c \cdot n_2^2 + c \cdot n^2$  bit operations. All values  $n$  from 1 to  $N$  are equally probable, so the average number of bit operations is equal to

$$c \cdot n_1^2 + c \cdot n_2^2 + c \cdot \frac{1}{N} \cdot \sum_{n=1}^N n^2. \quad (1)$$

For large  $N$ , the last sum in the expression (1) is approximately equal to  $N^3/3$ , so the expression (1) takes the form

$$c \cdot \left( n_1^2 + n_2^2 + \frac{N^2}{3} \right). \quad (2)$$

When  $n_1 \ll N$  and  $n_2 \ll N$ , the number of bit operations in the proposed setting is approximately equal to  $c \cdot (N^2/3)$  and is, thus, three times smaller than in the traditional setting.

**First conclusion.** *So indeed, the proposed setting decreases the number of bit operations by the factor of three.*

*Comment.* From the viewpoint of algorithmic complexity, making computations three times faster is not a big deal: most algorithms provide much more drastic speed-up; see, e.g., [1]. But let us recall that our goal is to extend the between-charges time for mobile computational devices. From this viewpoint, increasing the time-to-next-charging by a factor of three – e.g., from 8 hours to 24 hours – is a significant increase.

**A natural question: can we do even better?** We decreased the number of bit operations by a factor of three. A natural question is: is this the best we can do? Or can we decrease the number of bit operations even more?

To answer this question, let us consider the ideal situation when we know the exact precision  $n$  needed for each computational problem. In this ideal case, for each value  $n$ , we need exactly  $c \cdot n^2$  bit operations. Thus, the average number of bit operations is equal to

$$c \cdot \frac{1}{N} \cdot \sum_{n=1}^N n^2.$$

We already know that for large  $N$ , this expression is approximately equal to  $c \cdot (N^2/3)$  – which is exactly what our setting provides. Thus, we can make the following conclusion.

**Second conclusion.** *The proposed method is asymptotically optimal: it provides (asymptotically) the smallest possible number of bit operations – and thus, the smallest possible energy consumption that we can achieve without improving the algorithms.*

## Acknowledgments

This work was supported in part by the National Science Foundation grants 1623190 (A Model of Change for Preparing a New Generation for Professional Practice in Computer Science), HRD-1834620 and HRD-2034030 (CAHSI Includes), EAR-2225395 (Center for Collective Impact in Earthquake Science C-CIES), and by the AT&T Fellowship in Information Technology.

It was also supported by the program of the development of the Scientific-Educational Mathematical Center of Volga Federal District No. 075-02-2020-1478, and by a grant from the Hungarian National Research, Development and Innovation Office (NRDI).

The authors are thankful to Yuriy Kondratenko for his encouragement.

## References

- [1] Th. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*, MIT Press, Cambridge, Massachusetts, 2022.
- [2] R. Feynman, R. Leighton, and M. Sands, *The Feynman Lectures on Physics*, Addison Wesley, Boston, Massachusetts, 2005.
- [3] E. T. Jaynes and G. L. Bretthorst, *Probability Theory: The Logic of Science*, Cambridge University Press, Cambridge, UK, 2003.
- [4] C. Lauter and A. Volkova, “A framework for semi-automatic precision and accuracy analysis for fast and rigorous deep learning”, *Proceedings of the 2020 IEEE 27th Symposium on Computer Arithmetic (ARITH)*, Portland, Oregon, USA, June 7–10, 2020, pp. 103–110.
- [5] K. S. Thorne and R. D. Blandford, *Modern Classical Physics: Optics, Fluids, Plasmas, Elasticity, Relativity, and Statistical Physics*, Princeton University Press, Princeton, New Jersey, 2021.