



# Stochastic rounding and reduced-precision fixed-point arithmetic for solving neural ordinary differential equations

## Document Version

Final published version

[Link to publication record in Manchester Research Explorer](#)

## Citation for published version (APA):

Hopkins, M., Mikaitis, M., Lester, D., & Furber, S. (2020). Stochastic rounding and reduced-precision fixed-point arithmetic for solving neural ordinary differential equations. *Philosophical Transactions of the Royal Society A-Mathematical Physical and Engineering Sciences*, 1-22.

## Published in:

Philosophical Transactions of the Royal Society A-Mathematical Physical and Engineering Sciences

## Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

## General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

## Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact [uml.scholarlycommunications@manchester.ac.uk](mailto:uml.scholarlycommunications@manchester.ac.uk) providing relevant details, so we can investigate your claim.



Research



**Cite this article:** Hopkins M, Mikaitis M, Lester DR, Furber S. 2020 Stochastic rounding and reduced-precision fixed-point arithmetic for solving neural ordinary differential equations. *Phil. Trans. R. Soc. A* **378**: 20190052. <http://dx.doi.org/10.1098/rsta.2019.0052>

Accepted: 30 September 2019

One contribution of 15 to a discussion meeting issue 'Numerical algorithms for high-performance computational science'.

**Subject Areas:**

artificial intelligence, software, differential equations

**Keywords:**

fixed-point arithmetic, stochastic rounding, Izhikevich neuron model, ordinary differential equation, SpiNNaker, dither

**Authors for correspondence:**

Michael Hopkins

e-mail: [michael.hopkins@manchester.ac.uk](mailto:michael.hopkins@manchester.ac.uk)

Mantas Mikaitis

e-mail: [mantas.mikaitis@manchester.ac.uk](mailto:mantas.mikaitis@manchester.ac.uk)

<sup>†</sup>These authors contributed equally to this study.

# Stochastic rounding and reduced-precision fixed-point arithmetic for solving neural ordinary differential equations

Michael Hopkins<sup>†</sup>, Mantas Mikaitis<sup>†</sup>, Dave R. Lester and Steve Furber

APT research group, Department of Computer Science, The University of Manchester, Manchester, UK

 SF, 0000-0002-6524-3367

Although double-precision floating-point arithmetic currently dominates high-performance computing, there is increasing interest in smaller and simpler arithmetic types. The main reasons are potential improvements in energy efficiency and memory footprint and bandwidth. However, simply switching to lower-precision types typically results in increased numerical errors. We investigate approaches to improving the accuracy of reduced-precision fixed-point arithmetic types, using examples in an important domain for numerical computation in neuroscience: the solution of ordinary differential equations (ODEs). The Izhikevich neuron model is used to demonstrate that rounding has an important role in producing accurate spike timings from explicit ODE solution algorithms. In particular, fixed-point arithmetic with stochastic rounding consistently results in smaller errors compared to single-precision floating-point and fixed-point arithmetic with round-to-nearest across a range of neuron behaviours and ODE solvers. A computationally much cheaper alternative is also investigated, inspired by the concept of *dither* that is a widely understood mechanism for providing resolution below the least significant bit in digital signal processing. These results will have implications for the solution of ODEs in other subject areas, and should also be directly relevant to the huge range of practical problems that are represented by partial differential equations.

## 1. Introduction and motivation

Sixty-four-bit double-precision floating-point numbers are the accepted standard for numerical computing because they largely shield the user from concerns over numerical range, accuracy and precision. However, they come at a significant cost. Not only does the logic required to perform an arithmetic operation such as multiply-accumulate on this type consume an order of magnitude more silicon area and energy than the equivalent operation on a 32-bit integer or fixed-point type, but each number also requires double the memory storage space and double the memory bandwidth (and hence energy) each time it is accessed or stored. Today energy-efficiency has become a prime consideration in all areas of computing, from embedded systems through smartphones and data centres up to the next-generation exascale supercomputers. The pressure to pursue all avenues to improved energy efficiency is intense, causing all assumptions, including those regarding optimum arithmetic representations, to be questioned.

Nowhere is this pressure felt more strongly than in machine learning, where the recent explosion of interest in deep and convolutional nets has led to the development of highly effective systems that incur vast numbers of computational operations, but where the requirements for high-precision are limited. As a result, reduced precision types, such as fixed-point real types of various word lengths, are becoming increasingly popular in architectures developed specifically for machine learning. Where high throughput of arithmetic operations is required, accuracy is sacrificed by reducing the working numerical type from floating- to fixed-point and the word length from 64- or 32-bit to 16- or even 8-bit precision [1]. By reducing the word length, precision is compromised to gain the advantage of a smaller memory footprint and smaller hardware components and buses through which the data travel from memory to the arithmetic units. On the other hand, changing the numerical type from floating- to fixed-point, significantly reduces the range of representable values, increasing the potential for the under/overflow of the data types. In some tasks, these range issues can be dealt with safely by analysing the algorithm in various novel or problem-specific ways, whereas in other tasks it can be a very difficult to generalize effectively across all use cases [2]. There are other approaches apart from floating- and fixed-point arithmetics that are worth mentioning: *posit arithmetic* [3] is a completely new format proposed to replace floats and is based on the principles of *interval arithmetic* and *tapered arithmetic* [4] (dynamically sized exponent and significand fields which optimize the relative accuracy of the floating-point format in some specific range of real numbers rather than having the same relative accuracy across the whole range); *bfloat16*, with hardware support in recent Intel processors [5], is simply a single-precision floating-point type with the 16 bottom bits dropped for hardware and memory efficiency; *flexpoint* [6], an efficient combination of fixed- and floating-point also by Intel; and various approaches to transforming floating-point using, for example the *logarithmic number system* in which multiplication becomes addition [7].

A 32-bit integer adder uses  $9\times$  less energy [8] and  $30\times$  less area [9] than a single-precision floating-point adder. As well as reducing the precision and choosing a numerical type with a smaller area and energy footprint, mixed-precision arithmetic [10], stochastic arithmetic [11] and approximate arithmetic [12] have also been explored in the machine learning community with the goal of further reducing the energy and memory requirements of accelerators. Mixed-precision and stochastic techniques help to address precision loss when short word length numerical types are chosen for the inputs to and outputs from a hardware accelerator, and we address these in the present work. Mixed-precision arithmetic maintains intermediate results in formats different from the input and output data, whereas stochastic arithmetic works by using probabilistic rounding to balance out the errors in conversions from longer to shorter numerical types. These approaches

can also be combined. Approximate arithmetic (or more generally approximate computing) is a recent trend which applies the philosophy of adding some level of (application tolerant) error at the software level or inside the arithmetic circuits to reduce energy and area.

Interestingly, similar ideas have been explored in other areas, with one notable example being at the birth of the digital audio revolution where the concept of *dither* became an important contribution to understanding the ultimate low-level resolution of digital systems [13,14]. Practical and theoretical results show that resolution is available well below the least significant bit (LSB) if such a system is conceived and executed appropriately. Although the application here is different, the ideas are close enough to be worth mentioning and we show results with dither added at the input to an ODE solver. Similarly, an *approximate dithering adder*, which alternates between directions of error bounds to compensate for individual errors in the accumulation operation, has been reported [15].

Our main experiments are run on the spiking neural network (SNN) simulation architecture *SpiNNaker*. SpiNNaker is based on a digital 18-core chip designed primarily to simulate sparsely connected large-scale neural networks with weighted connections and spike-based signalling—a building block of large-scale digital neuromorphic systems that emphasizes flexibility, scaling and low energy use [16,17]. At the heart of SpiNNaker lie standard ARM968 CPUs. The ARM968 was chosen as it provides an energy-efficient processor core, addressing the need to minimize the power requirements of such a large-scale system. As the core supports only integer arithmetic, fixed-point arithmetic is used to handle real numbers unless the significant speed and energy penalties of software floating-point libraries can be tolerated. While many of the neural models that appear in the neuroscience literature can be simulated with adequate accuracy using fixed-point arithmetic and straightforward algorithms, other models pose challenges and require a more careful approach. In this paper, we build upon some previously published work [18] where the Izhikevich neuron model [19] was solved using fixed-point arithmetic, and show a number of improvements.

The main contributions of this paper are

- Improved rounding in fixed-point arithmetic: GCC implements rounding down in fixed-point arithmetic by default. This includes the rounding of constants in decimal to fixed-point conversion, as well as rounding in multiplication. Together these had a negative impact on accuracy in our earlier work [18]. In order to remedy this, we have implemented our own multiplication routines (including mixed-precision) with rounding control and used those instead of the GCC library routines (§4).
- The accuracy of ODE solvers for the Izhikevich neuron equations is assessed using different rounding routines applied to algorithms that use fixed-point types. This builds on earlier work where the accuracy of various fixed-point neural ODE solvers was investigated for this model [18], but where the role of rounding was not addressed (§5).
- Fixed-point ODE solvers are shown to be more robust with *stochastic rounding* than with other rounding algorithms, and are shown experimentally to be more accurate than single-precision floating-point ODE solvers (§5).
- 6 bits in the residual and random number are shown to be sufficient for a high-performance stochastic rounding algorithm when used in an ODE solver (§5c(iii)).
- The application of *dither* at the input to the ODE is shown to be a simpler and cheaper but, in many cases, effective alternative to full stochastic rounding (§6).

It should be noted that in this work when we talk about *error* we mean *arithmetic error*. This means the error introduced by imperfections in arithmetic calculations relative to an arithmetic reference, not the error introduced at a higher level by an algorithm relative to a better algorithm. We do not claim that a chosen algorithm is perfect, just that it is a realistic use case on which to test the different arithmetics that are of interest. This point is explained in more depth in §5a(i).

## 2. Background

### (a) Arithmetic types used in this study

There are eight arithmetic types under investigation in this study. Two of the floating-point types are standardized by the IEEE and widely available on almost every platform, and the other six are defined in the relatively recent ISO standard 18037 outlining C extensions for embedded systems [20] and are implemented by the GCC compiler:

- IEEE 64-bit double-precision floating-point: *double*
- IEEE 32-bit single-precision floating-point: *float*
- ISO 18037 s16.15 fixed-point: *accum*
- ISO 18037 u0.32 fixed-point: *unsigned long fract* or *ulf*
- ISO 18037 s0.31 fixed-point: *signed long fract* or *lf*
- ISO 18037 s8.7 fixed-point: *short accum*
- ISO 18037 u0.16 fixed-point: *unsigned fract* or *uf*
- ISO 18037 s0.15 fixed-point: *signed fract* or *f*

These types all come with the most common arithmetic operations implemented. The combination of these types and operations (rounded in various ways where appropriate) is the space that we explore in this paper.

### (b) Fixed-point arithmetic

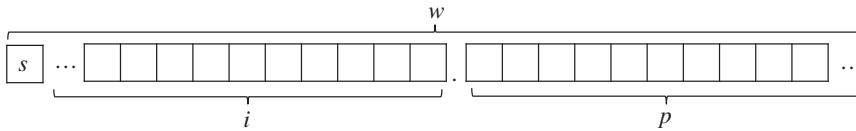
A **generalized fixed-point number** can be represented as shown in figure 1. We define a signed fixed-point number  $\langle s, i, p \rangle$  and an unsigned fixed-point number  $\langle u, i, p \rangle$  of word length  $w$  (which usually is 8, 16, 32 or 64) where  $i$  denotes the number of integer bits,  $p$  denotes the number of fractional bits and  $s$  is a binary value depending on whether a specific number is signed or unsigned (where  $u$  means that there is no sign bit at all and  $s$  means there is a sign bit and then we have to use 2's complement in interpreting the numbers).  $p$  also tells us the **precision** of the representation. Given a signed fixed-point number  $\langle s, i, p \rangle$ , the range of representable values is  $[-2^i, 2^i - 2^{-p}]$ . Whereas, given an unsigned fixed-point number  $\langle u, i, p \rangle$ , the range of representable values is  $[0, 2^i - 2^{-p}]$ . To measure the **accuracy** of a given fixed-point format (or more specifically some function that works in this format and we want to know how well, how *accurately* it performs in a given precision), we define *machine epsilon*  $\epsilon = 2^{-p}$ .  $\epsilon$  gives the smallest positive representable value in a given fixed-point format and therefore represents a *gap* or a *step* between any two neighbouring values in the representable range. Note that this gap is absolute across the representable range of values and is not relative to the *exponent* as in floating-point or similar representations. This requires us to consider only absolute error when measuring the accuracy of functions that are implemented in fixed-point representation. Accuracy is sometimes also measured in terms of LSB—a value represented by the LSB in a fixed-point word, which is the same as machine epsilon. Note that the maximum error of a fixed-point number, when round-to-nearest is used on conversion, is  $(\epsilon/2)$ .

Lastly, it is worth noting how to convert a general fixed-point number into a decimal number. Given a 2's complement fixed-point number of radix 2 (binary vector)  $\langle s, i, p \rangle$ :  $sI_{i-1}I_{i-2} \cdots I_0.F_1F_2 \cdots F_p$ , if the number is signed the decimal value is given by

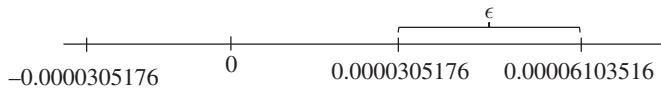
$$\text{value} = \sum_{k=0}^{i-1} I_k 2^k + \sum_{j=1}^p F_j 2^{-j} - s 2^i. \quad (2.1)$$

Otherwise, if the number is not signed (so bit  $s$  becomes integer bit  $I_i$ ) the decimal value is given by

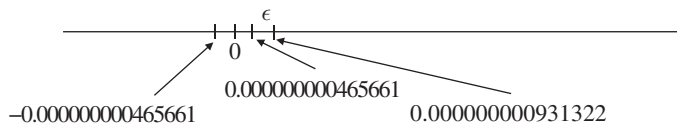
$$\text{value} = \sum_{k=0}^i I_k 2^k + \sum_{j=1}^p F_j 2^{-j}. \quad (2.2)$$



**Figure 1.** We define a fixed-point number that is made out of three parts: the most significant bit is a sign bit 's', an integer part with  $i$  bits and a fractional part with  $p$  bits. The *word length*  $w = i + p + 1$ .



**Figure 2.** Values close to zero in the *accum* representation.



**Figure 3.** Values close to zero in the *long fract* representation.

SpiNNaker software mostly uses two fixed-point formats available in the GCC implementation: *accum*, which is  $\langle s, 16, 15 \rangle$  and *long fract* which is  $\langle s, 0, 31 \rangle$ . The values close to 0 for each format are shown in figures 2 and 3. The *accum* type has a range of representable values of  $[-2^{16} = -65536, 2^{16} - 2^{-15} = 65535.99996948 \dots]$  with a gap between neighbouring values of  $\epsilon = 2^{-15} = 0.0000305175 \dots$ . The *long fract* type has a range of  $[-1, 1 - 2^{-31} = 0.9999999953433 \dots]$  with a gap of  $\epsilon = 2^{-31} = 0.0000000046566 \dots$  between neighbouring values. As can be seen from the values of machine epsilon, *long fract* is a more precise fixed-point representation. However, *long fract* has a very small range of representable values compared to *accum*. Which format should be used depends on the application requirements, such as the required precision and the bounds of all variables.

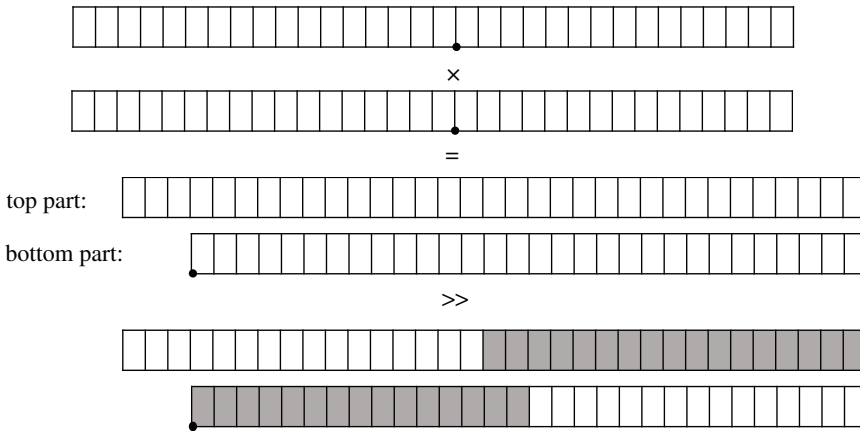
The three main **arithmetic operations** on fixed-point numbers can be defined as

- Addition:  $\langle s, i, p \rangle + \langle s, i, p \rangle = \langle s, i, p \rangle$ .
- Subtraction:  $\langle s, i, p \rangle - \langle s, i, p \rangle = \langle s, i, p \rangle$ .
- Multiplication:  $\langle s, i_a, p_a \rangle \times \langle s, i_b, p_b \rangle = \langle s, i_a + i_b, p_a + p_b \rangle$ .

Note that  $+$ ,  $-$  and  $\times$  denote integer operations available in most of processors' arithmetic logic units, including that of the ARM968. Therefore, for addition and subtraction, no special steps are required (compared to integer operation) when implementing and these operations are exact if there is no underflow or overflow. However, for multiplication, if the operands  $a$  and  $b$  have the word lengths  $w_a$  and  $w_b$ , then the result will have the word length of  $w_a + w_b - 1$ . Therefore, after the integer multiplication we have to shift the result right to obtain the result in the same fixed-point representation as the operands (the example in figure 4 shows this for the *accum* type). This shifting results in the **loss of precision** and therefore an appropriate rounding step can be done to minimize the error.

### (c) Rounding

Here, we will describe two known rounding approaches that we have used to increase the accuracy of the fixed-point multiplication: *round-to-nearest* (henceforth called *RN*) and *stochastic rounding* (henceforth called *SR*) [11,21]; the latter is named stochastic due to its use of random



**Figure 4.** Example showing the multiplication of two *accum* type variables. The shaded part is the 32-bit result that has to be extracted, discarding the bottom and top bits.

numbers. Given a real number  $x$ , and  $\langle s, i, p \rangle$ , an output fixed-point format, *round-to-nearest* is defined as

$$\text{RN}(x, \langle s, i, p \rangle) = \begin{cases} \lfloor x \rfloor, & \text{if } \lfloor x \rfloor \leq x < \lfloor x \rfloor + \frac{\epsilon}{2}, \\ \lfloor x \rfloor + \epsilon, & \text{if } \lfloor x \rfloor + \frac{\epsilon}{2} \leq x < \lfloor x \rfloor + \epsilon, \end{cases} \quad (2.3)$$

where  $\lfloor x \rfloor$  is the truncation operation (discarding a number of bottom bits and leaving  $p$  fractional bits) which returns a number in  $\langle s, i, p \rangle$  format less than or equal to  $x$ . Note that we chose to implement *rounding-up* on the tie  $x = \lfloor x \rfloor + (\epsilon/2)$  because this results in a simple rounding routine that requires checking only the *Most Significant Bit (MSB)* of the truncated part to make a decision to round up or down. Other tie breaking rules such as *round-to-even* can sometimes yield better results, but we prefer the cheaper tie-breaking rule in this work.

Stochastic rounding is similar, but instead of always rounding to the nearest number, the decision about which way to round is non-deterministic and the probability of rounding up is proportional to the residual. Given all the values as before and, additionally, given a random value  $P \in [0, 1)$ , drawn from a uniform random number generator, SR is defined (similarly as in [11]) as

$$\text{SR}(x, \langle s, i, p \rangle) = \begin{cases} \lfloor x \rfloor, & \text{if } P \geq \frac{x - \lfloor x \rfloor}{\epsilon}, \\ \lfloor x \rfloor + \epsilon, & \text{if } P < \frac{x - \lfloor x \rfloor}{\epsilon}. \end{cases} \quad (2.4)$$

Lastly, we will denote the cases without rounding as *rounding-down* (henceforth called *RD*)—this is achieved by truncating/shifting(right) a 2's complement binary number. Note that this results in round-down in 2's complement, while in floating-point arithmetic, which does not use 2's complement, bit truncation results in round towards zero.

#### (d) Floating-point arithmetic

The IEEE 754-2008 standard radix-2 normalized single-precision floating-point number with a sign bit  $S$ , exponent  $E$  and an integer significand  $T$ , is defined as ([22, p. 51], slightly modified)

$$(-1^S) \times 2^{E-127} \times (1 + T \cdot 2^{-23}), \quad (2.5)$$

whereas a double-precision number is defined as

$$(-1^S) \times 2^{E-1023} \times (1 + T \cdot 2^{-52}). \quad (2.6)$$



**Table 1.** Some features of floating-point (normalized) and 32-bit fixed-point numerical types.

property	double precision	single precision	<i>s16.15</i>	<i>u0.32</i>
accuracy	$2^{-52}$ (Rel.)	$2^{-23}$ (Rel.)	$2^{-15}$ (Abs.)	$2^{-32}$ (Abs.)
min (exact)	$2^{-1022}$	$2^{-126}$	$2^{-15}$	$2^{-32}$
min (approx.)	$2.225 \times 10^{-308}$	$1.175 \times 10^{-38}$	0.0000305	$2.32831 \times 10^{-10}$
max (exact)	$(2 - 2^{-52}) \times 2^{1023}$	$(2 - 2^{-23}) \times 2^{127}$	$2^{16} - 2^{-15}$	$1 - 2^{-32}$
max (approx.)	$1.798 \times 10^{308}$	$3.403 \times 10^{38}$	65535.999969	0.999...

Table 1 shows the minimum and maximum values of various floating- and fixed-point numerical types. It can be seen that fixed-point types have absolute accuracy denoted by the corresponding machine epsilon, which means that any two neighbouring numbers in the representable range have a fixed gap between them. On the other hand, floating-point types have accuracy relative to the exponent, which means that the gap between any two neighbouring numbers (or a real value of the LSB) is relative to the exponent that those numbers have, i.e. the corresponding machine epsilon multiplied by 2 to the power of the biased exponent. For example, the next number after 0.5 ( $E = 126$ ) in a single-precision floating-point number is  $0.5 + 2^{-23} \times 2^{-1}$ , while the next number after 1.0 ( $E = 127$ ) is  $1 + 2^{-23} \times 2^0$ . Owing to this, the accuracy of floating-point numbers is measured relative to the exponent.

Owing to these features of floating-point arithmetic, it depends on the application which types will give more accuracy overall. For example, if the application works with numbers around 1 only, *u0.32* is a more accurate type as it gives smaller steps between adjacent numbers ( $2^{-32}$ ) than does single-precision floating-point ( $2^{-23}$ ). On the other hand, if the application works with numbers that are as small as 0.001953125 ( $E = 117$ , which gives a gap between numbers of  $2^{-32}$ ), single-precision floating-point representation becomes as accurate as *u0.32* and increasingly more accurate as the numbers decrease beyond that point, eventually reaching the accuracy of  $2^{-126}$  (normalized) and  $2^{-126-23}$  (denormalized) near zero.

Therefore, it is very hard to show analytically what errors a complex algorithm will introduce at various steps of the computation for different numerical types. More generally, it is important to bear in mind that most of the literature on numerical analysis for computation implicitly assumes constant relative error (as in floating-point) as opposed to constant absolute error (as in fixed-point), and so many results on, for example, convergence and sensitivity of algorithms, need to be considered in this light. In this paper, we address these numerical issues experimentally using a set of numerical types with a chosen set of test cases.

### 3. Related work

There is already a body of work analysing reduced precision types and ways to alleviate numerical errors using stochastic rounding, primarily in machine learning research. In this section, we review some of the work that explores fixed-point ODE solvers and stochastic rounding applications in machine learning and neuromorphic hardware.

#### (a) Fixed-point ordinary differential equation solvers

The most recent work that explores fixed-point ODE solvers on SpiNNaker [23] was published in the middle of our current investigation and exposes some important issues with the default GCC *s16.15* fixed-point arithmetic when used with the Izhikevich neuron model. The authors tested the current sPyNNaker software framework [17] for simulating the Izhikevich neuron and then demonstrated a method for comparing the network statistics of two forward Euler solvers at small timesteps using a custom fixed-point type of *s8.23* for one of the constants. Comparing network behaviour is a valuable development in this area as is their observation about the need



for accurate constants in ODE solvers that we learned independently, as described later. However, we do have some comments on their methodology and note a few missing details in their study.

- (i) Although they mention the value of an absolute reference, they apparently do not use one in their quantitative comparison, and so they compare two algorithms neither of whose accuracy is established.
- (ii) The iterative use of forward Euler solvers with small time steps is far from optimal in terms of the three key measures of accuracy, stability and computational efficiency. Regarding the latter, we estimate that  $10\times$  as many operations are required compared to the RK2 solver that they compare against. Given this computational budget, a much more sophisticated solver could be used which would provide higher accuracy and stability than their chosen approach. The use of small-time steps also militates against the use of fixed-point types without re-scaling, because of quantization issues near the bottom of variable ranges that are inevitable with such small-time steps.
- (iii) The choice of the `s8.23` shifted type to improve the accuracy of the constant 0.04 seems odd when the `u0.32 long fract` type and matching mixed-format arithmetic operations can be used for any constants smaller than 1 and are fully supported by GCC, which could only result in a solver of equal or higher accuracy.
- (iv) Rounding is not explored as a possible improvement in the conversion of constants and arithmetic operations, and neither is the provision of an explicit constant for the nearest value that can be represented in the underlying arithmetic type.

In this study, we address all of these issues.

## (b) Stochastic rounding

Randomization of rounding can be traced back to the 1950s [24]. Furthermore, a similar idea was also explored in the context of the CESTAC method [25,26] where a program was run multiple times, making random perturbations of the last bit of the result of each elementary operation, and then taking a mean to compute the final result ([27, p. 486] and references therein).

Other studies investigate the effects of probabilistic rounding in backpropagation algorithms [21]. Three different applications are shown with varying degrees of precision in the internal calculations of the backpropagation algorithms. It is demonstrated that when fewer than 12 bits are used, training of the neural network starts to fail due to weight updates being too small to be captured by limited precision arithmetic, resulting in underflow in most of the updates. To alleviate this, the authors apply probabilistic rounding in some of the arithmetic operations inside the backpropagation algorithm and show that the neural network can then perform well for word widths as small as 4 bits. The authors conclude that with probabilistic rounding, the 12-bit version of their system performs as well as a single-precision floating-point version.

In a recent paper [11] about SR (and similarly in [28]), the authors demonstrate how the error resiliency of neural networks used in deep learning can be used to allow reduced precision arithmetic to be employed instead of the standard 32/64-bit types and operations. The authors describe a simple matrix multiplier array built out of a number of multiply–accumulate units that multiply two 16-bit results and accumulate the results in full-precision (implementing a dot product operation). Furthermore, the authors show that applying SR when converting the result of the dot product into lower 16-bit precision improves the neural network’s training and testing accuracy. The accuracy with a standard round-to-nearest routine is demonstrated to be very poor while stochastic rounding results are shown to be almost equal to the same neural network that uses single-precision floating-point. A very important result from this study is that 32-bit fixed-point types can be reduced to 16 bits and maintain neural network performance, provided that SR is applied. However, the authors did not report the effects of variation in the neural network results due to stochastic rounding—as applications using SR become stochastic themselves it is

important to report at least *mean* benchmark results and the *variation* across many trial runs as we have done in this study.

A recent paper from IBM [29] also explores the use of the 8-bit floating-point type with mixed-precision in various parts of the architecture and stochastic rounding. The authors demonstrate similar performance to the standard 32-bit float type in training neural networks.

In another paper [30], the effects of training recurrent neural networks—used in deep learning applications—on analogue resistive processing units (RPU) were studied. The authors aim to minimize the analogue hardware requirements by looking for a minimum number of bits that the *input arguments* to the analogue parts of the circuit can have. A baseline model with 5-bit input resolution is first demonstrated; it becomes significantly unstable (training error is reported) as network size is increased, compared to a simulation model with single-precision floating point. The authors then increase the input resolution to 7 bits and observe a much more regular development of the training error and with lower magnitude at the last training cycle. Finally, stochastic rounding is applied on these inputs and rounding them to 5 bits again makes the training error almost as stable as the 7-bit version, without the large training error observed in a 5-bit version that does not use stochastic rounding.

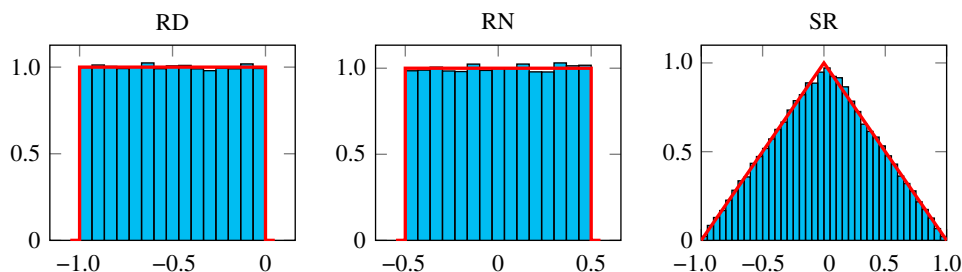
In neuromorphic computing, SR is used on the recently introduced Intel *Loihi* neuromorphic chip [31]. Here, SR is not applied to the ODE of the neuron model as it is in our study, but to the biological learning rule—spike-timing-dependent plasticity (STDP)—that defines how synapses change with spiking activity. The implementation of STDP usually requires spike history traces to be maintained for some number of timesteps. In *Loihi*, history traces (for each synapse) are held as 7-bit values to minimize memory cost and calculated as  $x[t] = \alpha \times x[t - 1] + \beta \times s[t]$ , where  $\alpha$  is a decay factor which defines how the trace decays back to 0 between spikes and  $\beta$  is a value contributed by each spike. The paper states that these traces are calculated using stochastic rounding.

## 4. SR implementation and testing with atomic multiplies

In order to establish the correctness of our rounding mechanisms for fixed-point types, we have carried out a set of tests to assess the distribution of errors from a simple atomic multiply with a wide range of randomly generated operands. There are five cases of interest (with five equivalent cases for 16-bit types):  $s16.15 \times s16.15$ ,  $s16.15 \times s0.31$ ,  $s16.15 \times u0.32$ ,  $u0.32 \times u0.32$ ,  $u0.32 \times s0.31$ . The first three of these return a saturated answer in  $s16.15$  and the last two in  $s0.31$ , with a sign used because some of the multiplication results in ODE solvers are preceded by a minus to invert the sign. Each multiplier has an option to do three types of rounding from RD, RN or SR when converting the result to a destination numerical type. In all cases, we take 50 000 random numbers distributed uniformly between a minimum and maximum value. For example, in the cases where both operands are  $s16.15$  the range is  $[-256, 256]$ , or  $[-16, 16]$  for the equivalent 16-bit types, to ensure that the result can be represented. In the other cases, the full ranges of the types can safely be used.

The operands are generated in pairs so that the numbers can be represented exactly in the relevant fixed-point type and then copied to the reference type IEEE *double* where the compiler translates the numerical value appropriately. We then multiply them using the default *double* operation and the other multiply operation that is under consideration. The fixed-point result is copied back to *double* and the difference between the results is stored and normalized to the LSB of the fixed point result type so that, e.g., if the result is an  $s16.15$ , the absolute error is multiplied by  $2^{15}$ . We then calculate summary statistics and plot each set of errors. We will call these *Bit Error Distributions* or *BEDs* and they allow the atomic multiplication errors of the different types to be compared directly, as well as to confirm the expected behaviour.

Figure 5 shows the BEDs of a  $s16.15 \times s16.15$  multiplier with the three different rounding modes. These plots and summary statistics calculated on the residuals confirm that in all cases the BEDs are as expected and we can have confidence in their implementation.



**Figure 5.** Histograms showing the bit error distribution of 50 000 random  $s16.15 \times s16.15$  operations with different rounding schemes. RD, RN and SR are round-down, round-to-nearest and stochastic rounding, respectively. (Online version in colour.)

### (a) Pseudo-random number generators

It is important to note that for SR to be competitive in terms of computation time a *very efficient source of high-quality pseudo-random numbers is critically important as in many cases it will be drawn on for every multiplication operation*. We have investigated three sources of uniformly distributed noise in this study, implemented as deterministic pseudo-random number generators (PRNGs) with varying levels of quality and execution speed.

The reference PRNG is a version of Marsaglia's KISS64 algorithm [32]. This has had several incarnations—we are using David Jones's implementation of what is called KISS99 here [33]. It is now the default PRNG on the SpiNNaker system with a long cycle length  $\approx 2^{123}$  and produces a high-quality random stream of 32-bit integers that satisfy the stringent TESTU01 BigCrush and Dieharder test suites [33,34]. Results have also been checked using faster PRNGs that fail these very challenging tests but which are considered to be of a good basic quality for non-critical applications. In reducing order of sophistication, these are a 33-bit maximum-cycle LFSR algorithm with taps at bits 33 and 20 implemented within the SpiNNaker base runtime library SARK [35] and a very simple Linear Congruential algorithm with a set up recommended by Knuth and Lewis, *ranq41* [36, p. 284]. No significant differences in either the mean or the standard deviation (SD) of the results were found in any case tested so far, which is encouraging in terms of showing that SR is insensitive to the choice of PRNG as long as it meets a certain minimum quality standard.

## 5. Ordinary differential equation solvers

The solution of ODEs is an important part of mathematical modelling in computational neuroscience, as well as in a wide variety of other scientific areas such as chemistry (e.g. process design), biology (e.g. drug delivery), physics, astronomy, meteorology, ecology and population modelling.

In computational neuroscience, we are primarily interested in solving ODEs for neuron behaviour though they can also be used for synapses and other more complex elements of the wider connected network such as gap junctions and neurotransmitter concentrations. Many of the neuron models that we are interested in are formally *hybrid systems* in that they exhibit both continuous and discrete behaviour. The continuous behaviour is defined by the time evolution of the internal variables describing the state of the neuron and the discrete behaviour is described by a threshold being crossed, a spike triggered followed by a reset of the internal state and then sometimes a refractory period set in motion whereby the neuron is temporarily unavailable for further state variable evolution.

A number of good references exist regarding the numerical methods required to solve ODEs when analytical solutions are not available, which is usually the case [37–40].

## (a) Izhikevich neuron

A sophisticated neuron model which provides a wide range of potential neuron behaviours is the one introduced by Eugene Izhikevich [19]. In our earlier work [18], we take a detailed look at this neuron model, popular because of its relative simplicity and ability to capture a number of different realistic neuron behaviours. Although we provide a new result in that paper which moves some way towards a closed-form solution, for realistic use cases we still need to solve the ODE using numerical methods. We aimed for a balance between the fidelity of the ODE solution with respect to a reference result and efficiency of computation using fixed-point arithmetic, with our primary aim being real-time operation on a fixed time grid (1 ms or 0.1 ms). The absolute reference used is that of the ODE solver in the Mathematica environment [41] which has many advantages including a wide range of methods that can be chosen adaptively during the solution, a sophisticated threshold-crossing algorithm, the ability to use arbitrary precision arithmetic and meaningful warning messages in the case of issues in the solution that contravene algorithmic assumptions. For the Izhikevich model, a well-implemented, explicit, variable-time-step solver such as the Cash–Karp variant of the Runge–Kutta fourth/fifth order will also provide a good reference as long as it implements an effective threshold crossing algorithm. For other models such as variants of the Hodgkin–Huxley model, it may be that more complex (e.g. implicit) methods will be required to get close to reference results. However, none of these will be efficient enough for real-time operation unless there is a significant reduction in the number of neurons that can be simulated, and so for the SpiNNaker toolchain we chose a variant of the fixed-time-step Runge–Kutta second-order solver that was optimized for fixed-point arithmetic to achieve an engineering balance between accuracy and execution speed.

### (i) Algorithmic error and arithmetic error

In our earlier work [18], we focused on *algorithmic error* i.e. how closely our chosen algorithm can match output from the reference implementation. This algorithmic error is created by the inability of a less accurate ODE solver method to match the reference result (assuming equivalent arithmetic). As the output in this kind of ODE model is fully described by the time evolution of the state variable(s), it is not easy to generate direct comparisons between these time-referenced vector quantities. There are few single number measures that are feasible to calculate while at the same time being useful to the computational neuroscience community. After some experiments, we decided on spike lead or lag, i.e. the timing of spikes relative to the reference. This relies on a relatively simple drift in one direction or the other for spike timings, but in all cases tested so far where the arithmetic error is relatively small it is reliable and easy to understand, measure and communicate. So we work with the lead or lag of spike times relative to the reference for a set of different inputs (either DC or exponentially falling to approximate a single synapse) and neuron types (three of the most commonly used Izhikevich specifications [19]: RS (regular spiking), CH (chattering) and FS (fast spiking)). We looked at different arithmetics and timesteps and their effect, but these were side issues compared to choosing and implementing an algorithm that balanced execution speed against fidelity to the reference using the fixed-point types defined by the ISO standard [20] at a 1 ms timestep.

As mentioned at the end of the introduction, it should be kept in mind that in this section, we are investigating a different matter: that of *arithmetic error*. Now we are in a position where we have chosen an algorithm for the above reasons and our new arithmetic reference is this algorithm calculated in IEEE double-precision arithmetic. So the purpose of this study is different: *how closely can other arithmetics come to this arithmetic reference?*

We provide results below on four different fixed-time-step algorithms in order to demonstrate some generality. In increasing order of complexity and execution time these are two second-order and one third-order fixed-timestep Runge–Kutta algorithms (e.g. [39,40]) and a variant of the Chan–Tsai algorithm [42]. All are implemented using the ESR (explicit solver reduction) approach described in [18] where the ODE definition and solver mechanism are combined and ‘unrolled’

into an equivalent algebraic solution which can then be manipulated algebraically in order to be optimized for speed and accuracy of implementation using the fixed-point types available. We refer to our earlier work [18] for other possibilities, and it is worth mentioning that we focused primarily there on a 1 ms timestep whereas for the current study we are using 0.1 ms. This is for a number of reasons: (1) it is becoming the new informal standard for accurate neural simulation; (2) accuracy is critically dependent upon the timestep chosen and so our solutions will now be significantly closer to the absolute algorithmic reference; and (3) relative differences in arithmetic should therefore be made more apparent.

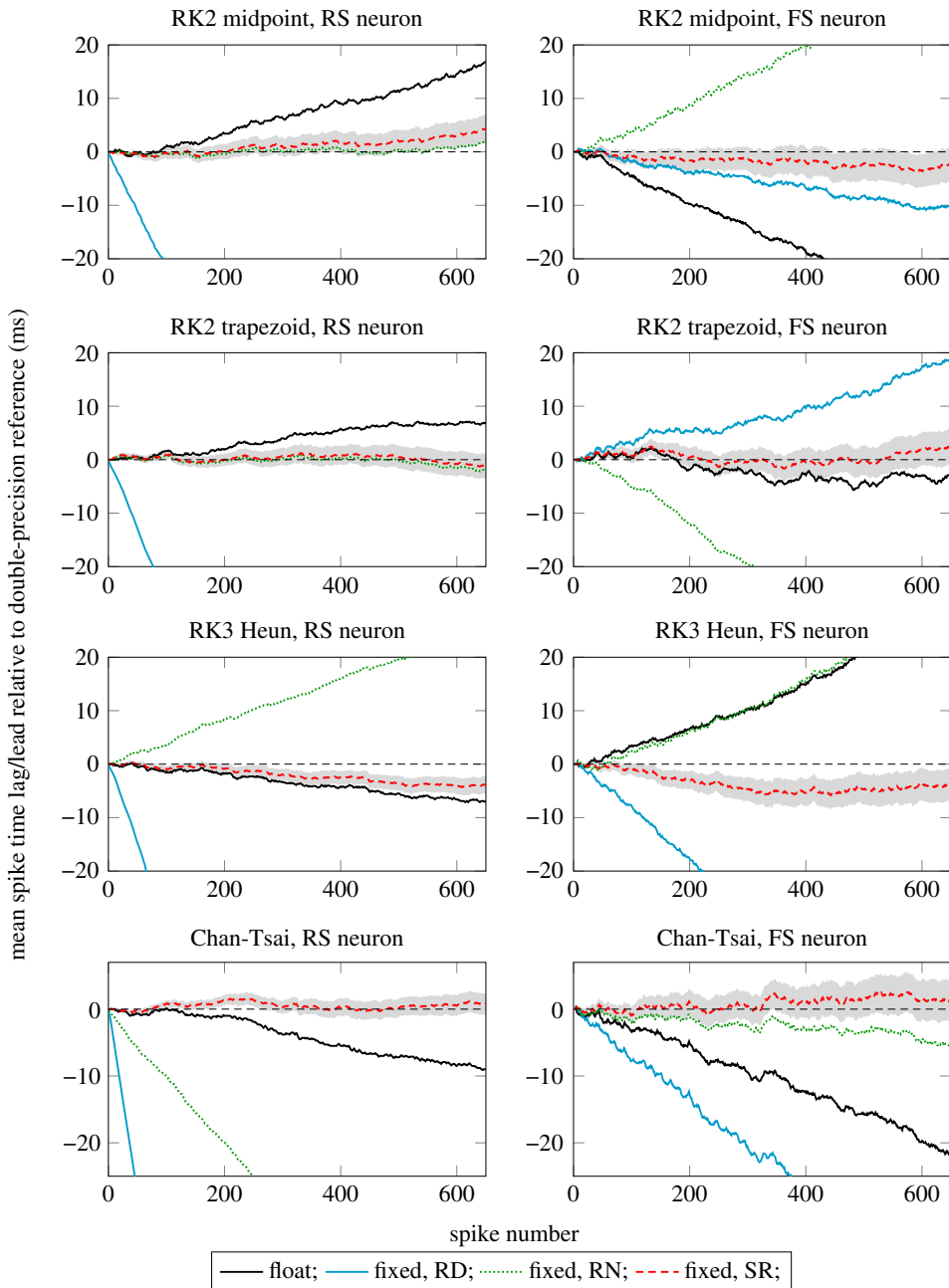
## (b) About constants

The constants that are used in the ODE solvers of the previously explored Izhikevich neuron model [19] cannot usually be represented exactly in a digital computer system. Originally, it was left to the compiler to round constants such as 0.04 and 0.1. However, as part of this work we have discovered that the GCC implementation of the fixed-point arithmetic *does not support rounding in decimal to fixed-point conversion* and therefore the specification of constants suffers from truncation error, which can be up to  $\epsilon$  for a given target fixed-point type, and for ODE solvers this error will in most cases accumulate. Additionally, we have found that *there is no rounding on most of the common arithmetic operations and conversions involving fixed-point numbers*. The pragma `FX_FULL_PRECISION` defined in the ISO standard [20] is not implemented in the GCC compiler version 6.3.1 that we have used for this work and therefore there is no way to turn on correct rounding. Experiments on GCC compiler v. 9.2.1 confirm that this is still the case.

Our first step was to specify constants in decimal exactly, correctly rounded to the nearest *s16.15* (e.g. explicitly stating 0.040008544921875 and 0.100006103515625 as constants instead of 0.04 and 0.1) which has reduced the maximum error in the decimal to *s16.15* conversion to  $(\epsilon/2) = (2^{-15}/2)$ . As a result, this has significantly reduced the spike lag previously reported [18], leading to an understanding that ODE solvers can be extremely sensitive to how the constants are represented. This was also noted in other work [23], and the solution that the authors took there was to add more bits in the fraction of the constants and add some scaling factors to return the final result in *s16.15*. Following from that, we have taken another step in this direction and have represented all of the constants smaller than 1 as *u0.32* instead of *s16.15*, which resulted in a maximum error of  $(2^{-32}/2)$ . The earlier work [23] used *s8.23* format for these constants, but we think that there is no downside to going all the way to *u0.32* format if the constants are below 1, and any arithmetic operations involving these constants can output in a different format if more dynamic range and signed values are required. In order to support this, we have developed libraries for mixed-format multiplication operations, where *s16.15* variables can be multiplied by *u0.32* variables returning an *s16.15* result (as described in detail in §4). A combination of more accurate representation of the constants, mixed-format multiplication and appropriate rounding has significantly reduced the error in the fixed-point solution of the Izhikevich neuron model from that previously reported [18] even before the 0.1 ms timestep is taken into account. Therefore, we would like to note that all the results below with *round-down* on the multipliers do not reproduce the results from the previous study [18] because these new results are generated with more precise constants, as well as some reordering of arithmetic operations in the originally reported ODE solvers in order to keep the intermediate values in high-precision fixed-point types for as long as possible.

## (c) Results

In this section, we present the results from four different ODE solvers with stochastic rounding added on each multiplier at the stage where the multiplier holds the answer in full precision and has to truncate it into the destination format, as described in §2. The experiments were run on the SpiNN3 board, which contains 4 SpiNNaker chips with ARM968 cores [16].



**Figure 6.** Spike lags of regular and fast-spiking Izhikevich neuron models for the DC input test at 0.1 ms timestep. Spike lags are computed against a double-precision floating-point reference in each case. The SR result is shown as the mean from 100 runs with the shaded area showing the SD. A negative value on the Y-axis indicates a lead, a positive value indicates a lag. (Online version in colour.)

### (i) Neuron spike timings

We show results here on the most challenging target problems. These are a constant DC input of  $\approx 4.775$  nA for the RS and FS neurons. The results for exponentially falling input and the CH neuron that were evaluated in our earlier work [18] are essentially perfect in terms of arithmetic error for all the possibilities; therefore we chose not to plot them. The results for four different solvers are shown in figure 6 with the exact spike lags of the 650th spike shown in table 2.



**Table 2.** Summary of ODE results: spike lags (ms) of the 650th spike in the DC input test. Positive indicates lag, negative - lead. Fixed, {RD/RN/SR} refer to fixed-point arithmetic with round-down, round-to-nearest and stochastic rounding, respectively.

solver	neuron type	float	fixed, RD	fixed, RN	fixed, SR (s.d.)
RK2 midpoint	RS	16.8	-131.4	1.9	4.3 (2.62)
	FS	-29.7	-10.0	33.7	-2.3 (3.16)
RK2 trapezoid	RS	6.9	-172.7	-2.1	-1.2 (2.30)
	FS	-3.2	18.7	-40.1	2.3 (3.33)
RK3 Heun	RS	-7.1	-206.9	26.0	-4.0 (1.59)
	FS	29.4	-53.6	31.4	-4.4 (3.10)
Chan-Tsai	RS	-9.0	-356.3	-67.9	0.8 (1.60)
	FS	-21.7	-44.6	-5.1	1.4 (3.10)

Because the SR results by definition follow a distribution, we need to show this in the plots. This has been done by running the ODE solver 100 times with a different random number stream and recording spike times on each run. We gather a distribution of spike times for each of the first 650 spikes and in figure 6 show the mean and SD of this distribution. A small number of checks with 1000 repeats show no significant differences in mean or SD from this case.

Clearly, the SR results are good in all of these cases compared to the alternatives, and in seven out of the eight cases are closest to the arithmetic reference after 650 spikes and look likely to continue this trend. All combinations of four different algorithms with disparate error sensitivities and two different ODE models are shown here, so this bodes well for the robustness of the approach. As a matter of interest, the RK2 trapezoid algorithm found to produce the most accurate solutions at 1 ms without correct rounding of constants or multiplications [18] continues to provide a good performance here in terms of arithmetic error, producing mean spike time errors of only -1.2 ms and +2.3 ms for the RS and FS neuron models after 69 and 165 seconds of simulation time, respectively. Perhaps unsurprisingly, one can see that the SR results look like a random walk.

## (ii) Evolution of the membrane potential $V$

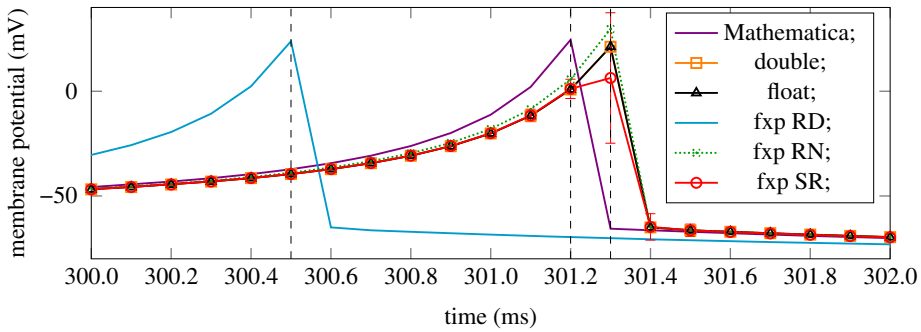
As an illustration of the imperfect evolution of the underlying state variables in the Izhikevich ODE, figure 7 shows the progression of the  $V$  state variable (the membrane potential of the neuron) after 300 ms for a variety of solver/arithmetic combinations and with a 0.1 ms timestep. One spike and reset event is shown in each case.

The absolute algorithmic reference given by Mathematica is shown in purple. The significant spike time lead given by RD is shown in the light blue line. The other results show a slight lag relative to the absolute reference, all clustered very close to the arithmetic reference in orange. The SR result shows the mean with SD error bars over 100 random results. The large SD at the spike event, and increased SD near it, are caused by a small number of realizations firing one or two timesteps early or late, inflating the SD sharply as the traces combine the threshold and reset voltages at these time steps. However, the mean behaviour tracks the arithmetic reference very closely.

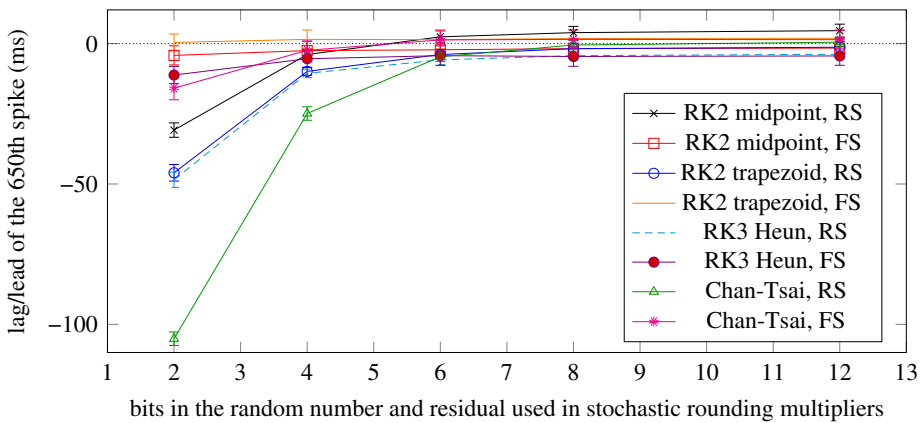
## (iii) The effects of reduced comparator precision in SR

In order to build an efficient SR hardware accelerator it is useful to consider how many bits are required in the rounding calculation. There are two ways to implement SR: the first is to build a comparator between the residual of the number to be rounded ( $(x - \lfloor x \rfloor)/(\epsilon)$ ) and the random number (as in equation (2.4)); the second approach that is strictly equivalent, implemented in an FPGA to optimize the utilization of the digital signal processing (DSP) units [11] is





**Figure 7.** The membrane potential of a neuron producing a third spike in the DC current test, using the RK2 Midpoint solver with various arithmetic and rounding formats. The temporal difference between the Mathematica and double-precision spike times is *algorithmic error* and the temporal difference between (for example) *fxp RD* and double is *arithmetic error*. (Online version in colour.)



**Figure 8.** Average spike lags of the 650th spike and the standard deviation from 100 runs. Four solvers with two different neuron types each are shown for varying numbers of bits in the stochastic rounding comparison step. All versions are compared to an equivalent reference solver implemented in double-precision floating-point arithmetic (the dotted straight line at  $y = 0$  ms). (Online version in colour.)

simply to add the random number to the input number and then truncate the fraction. This reverses equation (2.4) in such a way that the random numbers that are higher than or equal to  $1 - ((x - \lfloor x \rfloor) / (\epsilon))$  will produce a carry into the result bits (round-up) and the numbers that are lower will not produce a carry (which will result in round-down due to binary truncation). The first approach has a comparator of the remainder's/random number's width plus an adder after truncation of the remainder, whereas the second approach has only a single adder of the full word width. Whichever approach is used, the complexity of the adder/comparator can be reduced by defining a minimum number of bits required in the residual  $(x - \lfloor x \rfloor) / (\epsilon)$  and the random number.

In order to explore this problem in our ODE testbench, we ran a series of ODE solvers on two types of the Izhikevich neuron and measured the effects on the spike lag of different numbers of SR bits. The results are shown in figure 8 for the 650th spike. These results can be compared to the numbers in table 2 where all the available bits were used in SR. Note that the fixed-point multipliers in the ODE solvers can use 15 or 32 SR bits, depending on the types of the arguments. It is clear that as we decrease from 12 to 6 bits in the SR, the degradation in quality of the lead or lag relative to the arithmetic reference is negligible. However, when we decrease the number of bits from 6 to 4 or fewer, the regular spiking neuron starts to perform badly with a very high lag of the 650th spike. Interestingly, the fast-spiking neuron performs quite well even when only

**Table 3.** Summary of the ODE results with 16-bit fixed-point arithmetic compared to a double-precision reference (with the constants representable in 16 bits): spike lags (ms) of the 650th spike in the DC current test. Positive means lag, negative - lead. The test cases marked with a dash did not produce any spikes due to underflow in one of the internal calculations. Fixed, {RD/RN/SR} refer to fixed-point arithmetic with round-down, round-to-nearest and stochastic rounding, respectively.

solver	neuron type	float	fixed, RD	fixed, RN	fixed, SR (s.d.)
RK2 midpoint	RS	16.8	−21681.4	—	889.4 (58.82)
	FS	−29.7	−2754.5	686.4	676.7 (30.67)
RK2 trapezoid	RS	6.9	−22786.2	—	363.3 (57.65)
	FS	−4.6	−2391.2	892.8	516.7 (27.92)

2 SR bits are used. Given these two tests on the RS and FS neurons, we conclude that 4-bit SR is acceptable, with some degradation in quality of the neuron model, whereas the 6-bit version performs as well as the 12-bit or full remainder length SR version with 15 or 32 bits.

#### (iv) 16-bit arithmetic types

Using the same constant DC current test applied to a neuron, we evaluate 16-bit numerical types for this problem. It is a challenging task with such limited precision, but good results would improve memory and energy usage by a substantial further margin. The results with 16-bit ISO standard fixed-point types in the two second-order ODE solvers are shown in table 3. Because most of the variables are now held in *s8.7* numerical type, with 7 bits in the fractional part, even SR performs quite badly. However, it is clearly better than RD and RN, with RD performing terribly and RN being subject to underflow in one of the variables that causes updates to the main state variable to become 0 and therefore no spikes are produced. It seems that SR helps the ODE solver to recover from underflow and produce more reasonable answers with a certain amount of lag.

## 6. Ideas related to *Dither* and how these may be applicable

We mentioned in the introduction the idea of *dither* first applied in DSP for audio, image and video signals. The idea, in essence, is that as the input signal gets small enough to be close to the LSB of the digital system, some added noise helps to improve effective resolution of the system below the LSB. This gain is at the expense of a slight increase in broadband noise but, importantly for DSP applications, this noise is now effectively (or in certain cases *exactly*) uncorrelated with respect to the signal. This is in contrast to quantization errors which are highly correlated with the signal, and are considered to be objectionable and easy to identify in both audio and visual applications.

In our context, psychological considerations are difficult to justify. However, there are many other parallels between audio processing in particular and the solution of ODEs. Both are taking an input signal and using limited precision digital arithmetic operations in order to produce an output over a time sequence where any errors will be time-averaged. The ability of *dither* effectively to increase the resolution of the digital arithmetic operations was seen as an opportunity to discover if the benefits demonstrated earlier by SR on every multiply could be at least partially replicated by adding an appropriate quantity and distribution of random noise at the input to the ODE. In our cases, this is a DC current synaptic input, though any benefits would be generally applicable to the time-varying and conductance inputs that are also common in neural simulations.

A valuable text [43] contains much useful material about the precision and stability of finite precision computations. In particular, Chapter 8 describes the PRECISE toolbox which appears to use a similar idea to our application of *dither*. However, there are arguably two differences: (1) it is primarily a software technique for producing a number of sophisticated sensitivity analyses rather than a method for efficiently improving the accuracy of results and (2) averaging over many iterations in the solution of an ODE—and thereby cancelling out individual errors over time—is likely to be an important ingredient in the results described below.

Another analysis that helped lead to this idea was a consideration of the *backward error* of an algorithm [27, §1.5]. The most obvious consideration of error is that of *forward error* i.e. for a given input  $x$  and algorithm  $f(x)$  there is an output  $\hat{f}(x)$  which represents a realistic and imperfect algorithm that uses finite-precision arithmetic. The difference between  $f(x)$  and  $\hat{f}(x)$  is the forward error. However, it has been shown that in numerical work, the backward error is often a more useful measure. Imagine that you want to make  $\hat{f}(x)$  match the known ideal  $f(x)$  and all you are allowed to do is perturb  $x$  to  $\hat{x}$  so that theoretically  $\hat{f}(\hat{x}) \equiv f(x)$ . The difference between  $x$  and  $\hat{x}$  now describes the backward error and in many important real-world calculations bounds or distributions can be found for this value. It is important to consider how large this value is relative to the precision of your arithmetic type. If backward error is no greater than precision then the argument is that the error in a result from the algorithm is defined as much as, or more by, the precision of the input than it is by any imperfection in the algorithm itself.

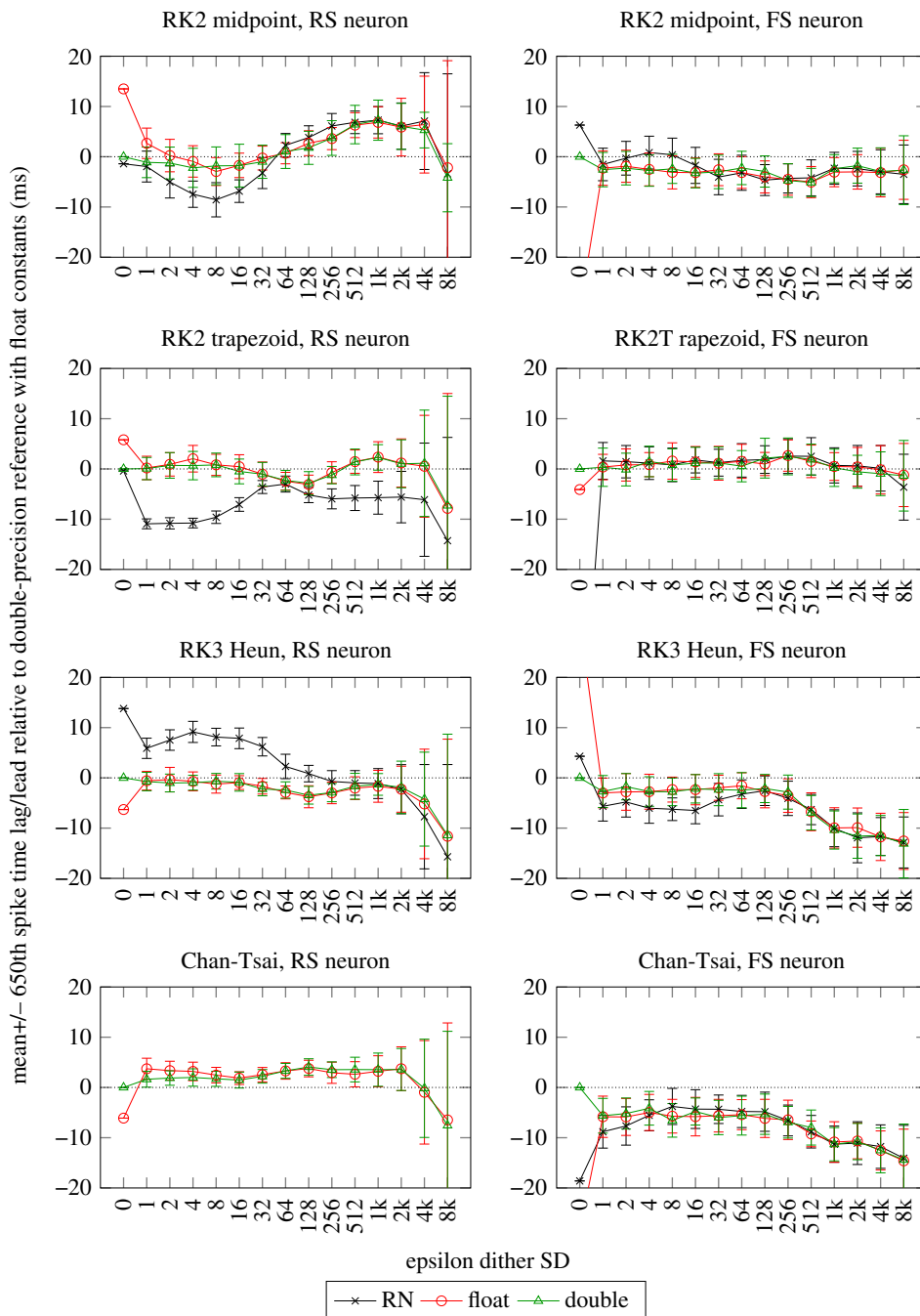
One can now relate this to the dither idea explored above. Now  $x$  can be seen as the input current to the neuron,  $f(x)$  as the ODE solution method which updates the state variable(s) at every timestep, and dither which adds random noise of some form and size to  $x$  becomes a probabilistic embodiment of the backward error of the algorithm. The hypothesis is that an appropriate dither signal introduced at the input can be found which is (1) no greater than the backward error, and (2) enough to ensure that stochasticity is induced in the arithmetic operations within the algorithm, thereby achieving indirectly a form of SR.

To identify appropriate values of dither, an empirical approach was used which adds Gaussian noise at the input scaled so that the SD ranges from 1 LSB for the s16.15 type (i.e.  $3.05e-5$ ) up to 10 000 LSBs. The DC input is 4.775 nA so that a dither value of 156 would give the input a Coefficient of Variation of 0.1%. Forty repeats of the eight neuron and solver combinations were then simulated. Unlike the earlier plots, 23-bit float constants were used in all cases which ensures that any offsets caused by differences in precision of the constants are corrected for and so any differences are now entirely due to the arithmetic differences caused by operations and how they are modified by the dither input.

Figure 9 shows the results with the epsilon multiplier on the X-axis with 0 meaning no dither and so equivalent to the reference results (except for the difference in constants for all except the *float* result). The Y-axis shows the mean difference in the 650th spike time from the *double* reference result over 40 runs with error bars showing the variation in this difference as a standard deviation. We have chosen a Y-axis scaling small enough to show subtle differences between the arithmetics, so in some cases very poor results are off the plot and will be described in the text.

The key points from these plots are the following:

- In all cases, *float* results are significantly improved by even the smallest amount of dither. This is particularly noticeable on the FS neuron but at this point there is no obvious explanation for this.
- In many cases, the smallest amount of dither also produces a significant improvement in the RN results. This is not the case for the RK2 solvers used with the RS neuron which perform well with 0 dither and small amounts of dither initially worsen the performance. The Chan-Tsai solver with the RS neuron has very high errors with RN and any amount of dither.
- The SD of the results is relatively consistent with large amounts of dither until a threshold, which is at approximately 1k for the RS neuron and approximately 4k for the FS neuron. The mean values drift away from the best results at approximately the same dither values. In some cases increasing input noise appears actually to reduce the spike lag SDs until they increase again at this threshold, e.g. both RK2 solvers with the RS neuron. This would seem counter-intuitive unless something other than random variation is relevant.
- Although not shown on the plots for clarity, SR tracks *double* better than RN and this is more apparent on the RS neuron. It is interesting to note in the full results how *float*, *double* and SR follow similar trends over dither value.



**Figure 9.** Effects of dither on DC input in the RS and FS neurons with four different solvers (single-precision floating-point constants). (Online version in colour.)

- There is an interesting apparent agreement between the results at approx 32–128 dither with all arithmetics (except RN in some cases) converging on the same result and perhaps the SD of the results also reducing slightly. This is more apparent on the RS neuron. There are a number of possible mechanisms for this effect—if it actually exists and is not just sampling variation—and it would be useful to understand them in terms of optimizing precision by choosing an appropriate dither value. There is an intriguing (but

speculative) possibility that results at these dither values are converging on a result with less algorithmic error than the *double* arithmetic reference which itself has finite precision. Further research is needed to confirm or refute these hypotheses.

- We have achieved similar results with stochastic rounding of all the constants involved in the ODE. The method was to choose one of the two versions of each constant on each ODE integration step stochastically, based on a more precise version of that constant. The results were approximately equivalent in second-order solvers.

Dither is in most cases going to be significantly simpler and cheaper than SR to implement, and for some applications this will be an important consideration. Runtime will be approximately  $N$  times faster if there are  $N$  operations within one iteration of the solver. For example, one of the ESR RK2 solvers investigated here uses 14 multiply operations. However, results so far are not as consistent or robust as SR. It appears to work better on the RK2 solvers. Our current hypothesis is that SR always induces the ideal uniform stochasticity for each multiply whereas in the case of dither the higher the order of the solver algorithm the more complex the processing of the input and hence the more distorted the Gaussian noise will become as it is increasingly transformed through successive arithmetic operations that make up the solver algorithm. An interim position where dither is added at various points within the algorithm may be of interest, but if added on too many operations it will become as expensive as the full SR solution or potentially more so because each Gaussian variate currently used for dithering is more difficult to generate than a uniform variate.

Recent work in large-scale climate simulation has explored something analogous to dither in our terms by introducing a noise source into the computational equations [44].

## 7. Discussion, further work and conclusion

In this paper, we addressed the numerical accuracy of ODE solvers, solving a well-known neuron model in fixed- and floating-point arithmetics. First, we identified that the constants in the Izhikevich neuron model should be specified explicitly by using the nearest representable number as the GCC fixed-point implementation does round-down in decimal to fixed-point conversion by default (this was also independently noticed by another study [23] but authors there chose to increase precision of the numerical format of the constants instead of rounding the constants to the nearest representable value as we did in this work). Next, we put all constants smaller than 1 into *unsigned long fract* types and developed mixed-format multiplications, instead of keeping everything in *accum*, to maximize the accuracy. This has not been done in any of the previous studies. We then identified that fixed-point multiplications are the main remaining source of arithmetic error and explored different rounding schemes. Round-to-nearest and stochastic rounding on multiplication results are shown to produce substantial accuracy improvements across four ODE solvers and two neuron types. Fixed-point with stochastic rounding was shown to perform better than fixed-point RN and float, and the mean behaviour is very close to double-precision ODE solvers in terms of spike times. We also found that simple PRNGs will often be good enough for SR to perform well. The minimum number of bits required in the random number in SR was found to be 6 across four different ODE solvers and two neuron types tested. In these cases, using more bits is unnecessary, and using fewer will cause the neuron timing to lead compared to the reference. Sixteen-bit arithmetic results were shown to have advantages with SR: although the absolute performance is quite poor (most likely due to overflows and underflows), 16-bit results with SR perform better than 16-bit RD (which largely results in spike time lead) and RN (which produces no spikes). Further work using scaled interim variables to ameliorate these issues is likely to provide further gains.

We investigated a method of adding noise to the inputs of the ODE on each integration step. Various levels of noise on the DC current input (dither) were shown to improve the accuracy of many of the solvers and is an alternative way, where high computational performance is a priority, to achieve some of the accuracy improvements that are shown with SR. Furthermore, we

found that stochastic rounding of all the constants, including the fixed timestep value, in each ODE integration step also improves spike timings in most of the test cases.

While the speed of ODE solvers is outside the scope of this paper, measurements in the test bench show that when SR is applied on multiplications, the speed is dictated by that of the random number generator. Preliminary numbers from RK2 Midpoint ODE solver performance benchmarks show an overhead of approximately 30% when RN is replaced with SR. Furthermore, SR is approximately  $2.6\times$  faster than software floating-point and approximately  $4.2\times$  faster than software double-precision in running a single ODE integration step.

The next generation SpiNNaker chip will have the KISS PRNG in hardware with new random numbers available in one clock cycle, so the overhead of SR compared to RN will be negligible. With regard to SpiNNaker-2, while it will have a single-precision floating-point unit, the results in this paper demonstrate that fixed-point with SR can be more accurate and should be considered instead of simply choosing floating point in a given application.

In the future, we aim to perform mathematical analysis to understand better why SR is causing less rounding error in ODEs as well as to understand in more detail how dither is improving accuracy in the majority of cases. For performance, we will explore how many random numbers are needed in each ODE integration step: do we require all multipliers to use a separate random number, or is it enough to use one per integration step? For SpiNNaker neuromorphic applications, we plan to build fast arithmetic libraries with SR and measure their overhead compared to default fixed-point arithmetic and various classical ways of rounding. Also, we plan to investigate the application of SR in neuron models other than Izhikevich's, and ways to solve them, the first example being LIF with current synapses which has a closed-form solution. Another direction is to investigate fixed-point arithmetic with SR in solving partial differential equations (PDEs) and other iterative algorithms (e.g. in Linear Algebra). Finally, it would be beneficial to investigate SR in reduced precision floating-point: *float16* and *bfloat16* numerical types which are becoming increasingly common in the machine learning and numerical algorithms communities, for large-scale projects such as weather simulation using PDEs [44,45].

Given all of these results, we predict that any reduced precision system solving ODEs and running other similar iterative algorithms will benefit from adding SR or noise on the inputs, but we would like to confirm this across a wider range of algorithms and problems. We would also like to point out that different arithmetics have different places where SR could be applied. For example, unlike fixed-point adders, floating-point adders and subtractors need to round when exponents do not match, and SR could be applied there after the addition has taken place (which would require preserving the bottom bits after matching the exponents). Similarly, in neural learning, where the computed changes to a weight are often smaller than the lowest representable value of that weight. To the best of our knowledge, the application of stochastic rounding in solving ODEs has not yet been investigated on any digital arithmetic, and these are the first results demonstrating substantial numerical error reduction in fixed-point arithmetic.

**Data accessibility.** The data and code used to generate the results presented in this paper are available from doi:10.17632/wxx2mf6n5n.1.

**Authors' contributions.** M.M. found the issue with rounding in GCC, suggested using SR in the Izhikevich neuron ODE solver and implemented fixed-point libraries with SR. M.H. implemented PRNGs, fixed-point ODE solvers and a test suite for the Izhikevich ODE. M.M. and M.H. extended the original ODE test suite, designed experiments and analysed the data. M.H. identified arithmetic error as the focal point and carried out the probabilistic analysis of the general SR case. M.H. and S.B.F. introduced the idea of dither. S.B.F. and D.R.L. provided supervision. All authors wrote the manuscript and approved the final submission.

**Competing interests.** S.B.F. is a founder, director and shareholder of Cogniscience Ltd, which owns SpiNNaker IP. M.H. and D.R.L. are shareholders of Cogniscience Ltd.

**Funding.** The design and construction of the SpiNNaker machine was supported by EPSRC (the UK Engineering and Physical Sciences Research Council) under grant nos. EP/D07908X/1 and EP/G015740/1, in collaboration with the universities of Southampton, Cambridge and Sheffield and with industry partners ARM Ltd, Silistix Ltd and Thales. Ongoing development of the software, including the work reported here, is supported by the EU ICT Flagship Human Brain Project no. (H2020 785907), in collaboration with many



university and industry partners across the EU and beyond. M.M. is funded by a Kilburn studentship at the School of Computer Science.

**Acknowledgements.** Thanks to Nick Higham for a valuable discussion about backward error analysis, Jim Garside for discussions about the number of bits required in the random numbers in stochastic rounding and to the reviewers for useful comments and references (particularly the work of Prof. Chaitin-Chatel in [43]) which have improved the paper.

## References

1. Jouppi NP *et al.* 2017 In-datacenter performance analysis of a tensor processing unit. In *Proc. of the 44th Annual Int. Symp. on Computer Architecture, ISCA '17*, pp. 1–12. New York, NY: ACM. Available from: <http://doi.acm.org/10.1145/3079856.3080246>.
2. Kabi B, Sahadevan AS, Pradhan T. 2017 An overflow free fixed-point eigenvalue decomposition algorithm: case study of dimensionality reduction in hyperspectral images. In *Conf. On Design And Architectures For Signal And Image Processing (DASIP). Dresden, Germany, 27–29 September*. Piscataway, NJ: IEEE. Available from: <http://dasip2017.esit.rub.de/program.html>.
3. Gustafson J, Yonemoto I. 2017 Beating floating point at its own game: posit arithmetic. *Supercomput. Front. Innov. Int. J.* **4**, 71–86. (doi:10.14529/jsf170206)
4. Morris R. 1971 Tapered floating point: a new floating-point representation. *IEEE Trans. Comput.* **C-20**, 1578–1579. (doi:10.1109/T-C.1971.223174)
5. Intel. BFLOAT16—Hardware Numerics Definition; 2018. Online: <https://software.intel.com/sites/default/files/managed/40/8b/bf16-hardware-numerics-definition-white-paper.pdf>.
6. Köster U *et al.* 2017 Flexpoint: an adaptive numerical format for efficient training of deep neural networks. In *Advances in neural information processing systems*, pp. 1742–1752. Available from: <https://arxiv.org/abs/1811.01721>.
7. Johnson J. 2018 Rethinking floating point for deep learning. (<http://arxiv.org/abs/quant-ph/181101721>).
8. Horowitz M. 2014 1.1 Computing's energy problem (and what we can do about it). In *2014 IEEE Int. Solid-State Circuits Conf. Digest of Technical Papers (ISSCC)*, pp. 10–14.
9. Dally W. 2016 High-Performance Hardware for Machine Learning. Online: [https://berkeley-deep-learning.github.io/cs294-dl-fl16/slides/DL\\_HW\\_Berkeley\\_0916.pdf](https://berkeley-deep-learning.github.io/cs294-dl-fl16/slides/DL_HW_Berkeley_0916.pdf).
10. Micikevicius P *et al.* 2018 Mixed precision training. In *Proc. of the 6th Int. Conf. on Learning Representations. ICLR'18, Vancouver, Canada, 30 April*. La Jolla, CA: ICLR.
11. Gupta S, Agrawal A, Gopalakrishnan K, Narayanan P. 2015 Deep learning with limited numerical precision. In *Proc. of the 32nd Int. Conf. on Int. Conf. on Machine Learning - vol. 37. ICML'15. JMLR.org*. pp. 1737–1746. Available from: <http://dl.acm.org/citation.cfm?id=3045118.3045303>.
12. Han J, Orshansky M. 2013 Approximate computing: an emerging paradigm for energy-efficient design. In *2013 18th IEEE European Test Symp. (ETS)*, pp. 1–6.
13. Vanderkooy J, Lipshitz SP. 1987 Dither in digital audio. *J. Audio Eng. Soc.* **35**, 966–975.
14. Vanderkooy J, Lipshitz SP. 1989 Digital dither: processing with resolution far below the least significant bit. In *Audio Engineering Society Conf. 7th Int. Conf. Audio in Digital Times. Toronto, Canada, 1 May*. New York, NY: Audio Engineering Society. Available from: <http://www.aes.org/e-lib/browse.cfm?elib=5482>.
15. Miao J, He K, Gerstlauer A, Orshansky M. 2012 Modeling and synthesis of quality-energy optimal approximate adders. In *2012 IEEE/ACM Int. Conf. on Computer-Aided Design (ICCAD)*, pp. 728–735. San Jose, CA, 5–8 November. Piscataway, NJ: IEEE.
16. Furber SB, Galluppi F, Temple S, Plana LA. 2014 The SpiNNaker project. *Proc. IEEE* **102**, 652–665. (doi:10.1109/JPROC.2014.2304638)
17. Rhodes O *et al.* 2018 sPyNNaker: a software package for running PyNN simulations on SpiNNaker. *Front. Neurosci.* **12**, 816. (doi:10.3389/fnins.2018.00816)
18. Hopkins M, Furber S. 2015 Accuracy and efficiency in fixed-point neural ODE solvers. *Neural Comput.* **27**, 2148–2182. (doi:10.1162/NECO\_a\_00772)
19. Izhikevich EM. 2003 Simple model of spiking neurons. *IEEE Trans. Neural Netw.* **14**, 1569–1572. (doi:10.1109/TNN.2003.820440)



20. ISO/IEC. 2008 ISO/IEC, Programming languages - C - Extensions to support embedded processors, ISO/IEC TR 18037:2008. ISO/IEC JTC 1/SC 22; Available from: <https://www.iso.org/standard/51126.html>.
21. Höhfeld M, Fahlman SE. 1992 Probabilistic rounding in neural network learning with limited precision. *Neurocomputing* **4**, 291–299. (doi:10.1016/0925-2312(92)90014-G)
22. Muller JM, Brunie N, de Dinechin F, Jeannerod CP, Joldes M, Lefèvre V, Melquiond G, Revol N, Torres S. 2018 *Handbook of floating-point arithmetic*, 2nd edn. ACM G.1.0; G.1.2; G.4; B.2.0; B.2.4; F.2.1. Basel, Switzerland: Birkhäuser.
23. Trensck G, Gutzen R, Blundell I, Denker M, Morrison A. 2018 Rigorous neural network simulations: a model substantiation methodology for increasing the correctness of simulation results in the absence of experimental validation data. *Front. Neuroinformatics* **12**, 81. (doi:10.3389/fninf.2018.00081)
24. Forsythe G. 1959 Reprint of a note on rounding-off errors. *SIAM Rev.* **1**, 66–67. (doi:10.1137/1001011)
25. Scott NS, Jézéquel F, Denis C, Chesneaux JM. 2007 Numerical ‘Health Check’ for scientific codes: the CADNA approach. *Comput. Phys. Commun.* **176**, 507–521. (doi:10.1016/j.cpc.2007.01.005)
26. Alt R, Vignes J. 2008 In *Stochastic Arithmetic as a Model of Granular Computing* (eds W Pedrycz, A Skowron, V Kreinovich), pp. 33–54. New York, NY: Wiley.
27. Higham NJ. 2002 *Accuracy and stability of numerical algorithms*, 2nd edn. Philadelphia, PA: Society for Industrial and Applied Mathematics.
28. Müller LK, Indiveri G. 2015 Rounding methods for neural networks with low resolution synaptic weights. (<https://arxiv.org/abs/1504.05767>)
29. Wang N, Choi J, Brand D, Chen CY, Gopalakrishnan K. 2018 Training deep neural networks with 8-bit floating point numbers. In *NeurIPS*. (<https://arxiv.org/abs/1812.08011>)
30. Gokmen T, Rasch MJ, Haensch W. 2018 Training LSTM networks with resistive cross-point devices. *Front. Neurosci.* **12**, 745. (doi:10.3389/fnins.2018.00745)
31. Davies M. 2018 Loihi: a neuromorphic manycore processor with on-chip learning. *IEEE Micro* **38**, 82–99. (doi:10.1109/MM.2018.112130359)
32. Marsaglia G, Zaman A. 1993 The KISS generator. Department of Statistics, Florida State University, Tallahassee, FL, USA.
33. L’Ecuyer P, Simard R. 2007 TestU01: a C library for empirical testing of random number generators. *ACM Trans. Math. Softw.* **33**, 22. (doi:10.1145/1268776.1268777)
34. Brown R. dieharder. Available from: <http://webhome.phy.duke.edu/~rgb/General/dieharder.php>.
35. Golomb SW. 1981 *Shift register sequences*. Laguna Hills, CA, USA: Aegean Park Press.
36. Press WH, Teukolsky SA, Vetterling WT, Flannery BP. 1992 *Numerical recipes in C*, 2nd edn. New York, NY: Cambridge University Press.
37. Hall G, Watt JM (eds). 1976 *Modern numerical methods for ordinary differential equations*. Oxford, UK: Clarendon Press.
38. Lambert JD. 1991 *Numerical methods for ordinary differential systems: the initial value problem*. New York, NY: John Wiley.
39. Dormand JR. 1996 *Numerical methods for differential equations*. Boca Raton, FL: CRC Press.
40. Butcher JC. 2003 *Numerical methods for ordinary differential equations*. New York, NY: John Wiley.
41. Wolfram Research. 2014 Mathematica. Version 10.0 ed. Champaign, IL: Wolfram Research, Inc.
42. Chan RPK, Tsai AYJ. 2010 On explicit two-derivative Runge–Kutta methods. *Numer. Algorithms* **53**, 171–194. (doi:10.1007/s11075-009-9349-1)
43. Chaitin-Chatelin F, Frayssé V. 1996 *Lectures on finite precision computations*. Philadelphia, PA: SIAM.
44. Palmer T. 2015 Build imprecise supercomputers. *Nature* **526**, 32–33. (doi:10.1038/526032a)
45. Dawson A, Düben PD, MacLeod DA, Palmer TN. 2018 Reliable low precision simulations in land surface models. *Clim. Dyn.* **51**, 2657–2666. (doi:10.1007/s00382-017-4034-x)