



## UvA-DARE (Digital Academic Repository)

### Periodic orbits in chaotic systems simulated at low precision

Klöwer, M.; Coveney, P.V.; Paxton, E.A.; Palmer, T.N.

**DOI**

[10.1038/s41598-023-37004-4](https://doi.org/10.1038/s41598-023-37004-4)

**Publication date**

2023

**Document Version**

Final published version

**Published in**

Scientific Reports

**License**

CC BY

[Link to publication](#)

**Citation for published version (APA):**

Klöwer, M., Coveney, P. V., Paxton, E. A., & Palmer, T. N. (2023). Periodic orbits in chaotic systems simulated at low precision. *Scientific Reports*, 13, Article 11410. <https://doi.org/10.1038/s41598-023-37004-4>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.



# OPEN Periodic orbits in chaotic systems simulated at low precision

Milan Klöwer<sup>1,2</sup>✉, Peter V. Coveney<sup>3,4,5</sup>, E. Adam Paxton<sup>1</sup> & Tim N. Palmer<sup>1</sup>

Non-periodic solutions are an essential property of chaotic dynamical systems. Simulations with deterministic finite-precision numbers, however, always yield orbits that are eventually periodic. With 64-bit double-precision floating-point numbers such periodic orbits are typically negligible due to very long periods. The emerging trend to accelerate simulations with low-precision numbers, such as 16-bit half-precision floats, raises questions on the fidelity of such simulations of chaotic systems. Here, we revisit the 1-variable logistic map and the generalised Bernoulli map with various number formats and precisions: floats, posits and logarithmic fixed-point. Simulations are improved with higher precision but stochastic rounding prevents periodic orbits even at low precision. For larger systems the performance gain from low-precision simulations is often reinvested in higher resolution or complexity, increasing the number of variables. In the Lorenz 1996 system, the period lengths of orbits increase exponentially with the number of variables. Moreover, invariant measures are better approximated with an increased number of variables than with increased precision. Extrapolating to large simulations of natural systems, such as million-variable climate models, periodic orbit lengths are far beyond reach of present-day computers. Such orbits are therefore not expected to be problematic compared to high-precision simulations but the deviation of both from the continuum solution remains unclear.

Many natural systems exhibit chaotic dynamics. The chaos in weather prevents reliable forecasts beyond one or two weeks<sup>1,2</sup>. The turbulent flow of air or water around vehicles requires complex numerical simulations to optimise the drag<sup>3,4</sup>. Similarly, chaos is present in models of many-body problems from astrophysics<sup>5,6</sup>, classical molecular dynamics and chemical reaction networks<sup>7</sup>, and plasma in fusion reactors<sup>8,9</sup>. As chaotic dynamics often prevent analytical solutions, numerical simulations with finite-precision floating-point numbers<sup>10</sup> are used to approximate a system's state and predict its future.

However, simulating a deterministic, yet chaotic system with deterministic finite-precision numbers always results in closed periodic orbits due to a finite set of possible states<sup>11</sup>. One may think of these orbits as *bitwise* periodic, in the sense that they eventually return to a state in which every bit of every independent variable is identical to the initial state. Eternal periodicity in the simulation of chaotic systems violates a fundamental property of chaos, but periods are very long with the high precision of 64-bit floating-point numbers (Float64).

Unstable periodic orbits are the *skeleton of chaos*<sup>12</sup> and have been intensively studied to better understand the dynamical properties of chaotic systems<sup>13–15</sup>. Chaotic trajectories follow a given periodic orbit in its vicinity but eventually diverge due to the orbit's instability and approach another orbit until diverging again<sup>16</sup>. The spectrum of the periodic orbits is a decomposition of the attractor<sup>17</sup>. It contains infinitely many countable orbits: the more orbits the longer the period with no upper bound on the period length<sup>16,18</sup>. The spectrum is truncated with deterministic finite-precision to a finite set of orbits bounded by a maximum period length. How such a truncation degrades the simulated dynamics is generally unclear. While very low-precision arithmetic truncates chaotic attractors to fixed points or short loops, the degradation can also be less obvious: it may be either substantial or to the point of being negligible.

Computed dynamics have errors relative to analytical solutions. Model errors arise from the difference between the mathematical equations and the natural systems they represent, including unresolved processes and heuristic parameters. Errors in the initial or boundary conditions are a result of imperfect observations being assimilated into the numerical model<sup>19</sup>. Discretization errors occur when a continuous system is discretized into a number of variables that is often limited by available computational resources<sup>20</sup>. In addition, there are rounding errors as a result of using finite-precision numbers to approximate real numbers<sup>21</sup>. In a chaotic system, errors

<sup>1</sup>Atmospheric, Oceanic and Planetary Physics, University of Oxford, Oxford, UK. <sup>2</sup>Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>3</sup>Centre for Computational Science, University College London, London, UK. <sup>4</sup>Advanced Research Computing Centre, University College London, London, UK. <sup>5</sup>Informatics Institute, University of Amsterdam, Amsterdam, The Netherlands. ✉email: milank@mit.edu

grow exponentially with time, so the largest source of error masks smaller ones. Due to often negligible rounding errors, many numerical simulations are currently transitioning to low-precision calculations in exchange for computational performance<sup>22–24</sup>. The performance gain is then reinvested into a higher resolution or complexity with more independent variables, typically with the intention of increasing the model's accuracy.

16-bit low-precision computations are increasingly supported on modern processors, such as graphics processing units<sup>25</sup> (GPU), tensor processing units<sup>26</sup> (TPU) and also conventional central processing units<sup>27,28</sup> (CPU). While the standard and only widely available number format are floats, several alternatives have been proposed: posits<sup>29</sup>, logarithmic fixed-point numbers<sup>30</sup> (logfix), and floats with stochastic rounding<sup>31–33</sup>. Currently lacking in hardware support, these number formats are first emulated in software for precision tests. The comparison across formats provides a better understanding as to how the numerical precision affects the simulated dynamics.

Here, we compare the periodic orbits and invariant measures as the properties of three chaotic dynamical systems when simulated with different binary number formats and rounding modes and at various levels of precision. The number formats are briefly described in section "Periodic orbits and number formats" along with our methodology for finding periodic orbits. In section "The logistic map with various number formats" we analyse the bifurcation of the logistic map as simulated with different number formats and rounding modes. In section "Revisiting the generalised Bernoulli map", we revisit the simulation of the generalised Bernoulli map to analyse numerical precision in a system where the analytical invariant measure is known. In section "Orbits in the Lorenz 1996 system", we turn to the Lorenz 1996 system to investigate the periodic orbit spectrum with an increasing number of variables. Section "Conclusions" summarises the results, section "Methods" provides further details on the methodology.

## Periodic orbits and number formats

The state of a deterministic dynamical system is entirely determined by  $X = (x_1, x_2, \dots, x_N)$ , the vector of all its  $N$  prognostic variables in a given finite-precision number format at a given time step  $t$ . We define a periodic orbit in a simulation when the state vector  $X_{t_0}$  at time step  $t_0$  exactly reoccurs at a later time step  $t_1 > t_0$ ,

$$X_{t_0} = X_{t_1}. \quad (1)$$

Equality is hereby required within the considered precision of the number format. Equivalently,  $X_{t_0}$  and  $X_{t_1}$  are bitwise identical. An exception occurs for floats where  $-0 = 0$ , which arithmetically does not impact on the dynamical system. This bitwise periodicity is in contrast to other studies investigating quasi-periodic orbits<sup>34,35</sup>, which require  $X_{t_0}, X_{t_1}$  to be close, but not bitwise identical. The periodic orbits here are found in long simulations when Eq. (1) holds and are very sensitive to the choice of the number format and numerical precision. This is distinct from unstable periodic orbits, numerically found via an iterative Newton method when Eq. (1) holds up to a numerical error<sup>36</sup> that is larger than the precision of the number format.

The invariant measure of a chaotic system describes its attractor independent of the initial conditions. A deterministic chaotic system simulated with deterministic finite-precision arithmetic will converge to one of the periodic orbits for any initial condition. Based on all periodic orbits we can compute the invariant measure through a weighted average by the orbits' respective basins of attraction, i.e. the fraction of initial conditions that end up on a given periodic orbit. To find *all* periodic orbits in a deterministic dynamical system, all possible initial conditions have to be integrated and checked for periodicity. For a 1-variable system with  $X \in [0, 1)$  (such as the Bernoulli map, see section "Revisiting the generalised Bernoulli map") simulated with Float32 there are 1,065,353,216 unique initial conditions. For any larger system or higher precision number format it becomes computationally virtually impossible to consider all initial conditions. We therefore use a Monte Carlo-based random sampling of the initial conditions to find a subset of all orbits. An improved pseudo-random number generator for floats is developed for this purpose (see Methods). The orbits found are expected to be those with the largest basins of attraction, and a robust estimate of their size is obtained for a sufficiently large sample of initial conditions. This procedure is explained in the section Methods.

Floating-point numbers are standardised following IEEE<sup>10,37</sup>. Other binary number formats are briefly introduced. Posits<sup>29</sup> have a slightly higher precision within the powers of 2 around  $\pm 1$ , yet a wide dynamic range at the cost of a gradually lower precision away from  $\pm 1$ . While posits have been proposed as a drop-in replacement for floats, they currently lack widely available hardware support. For more details on posits see refs.<sup>29,38–40</sup>. Logarithmic fixed-point numbers have received little attention apart from research implementations on custom hardware<sup>30</sup> and are therefore described with our design choices in more detail.

**Logarithmic fixed-point numbers.** The logarithmic fixed-point number (logfix) format LogFixPoint16 is here defined with a similar range-precision trade-off as Float16 with a sign bit  $s$ ,  $n_i = 5$  signed integer bits in the exponent and  $n_f = 10$  fraction bits. A logfix number  $x$  is of the form

$$x = (-1)^s \cdot 2^k \quad (2)$$

The exponent  $k = i + f$  is equivalent to a fixed-point number with  $n_f$  binary digits accuracy. The signed integer  $i$  is in  $[-2^{n_i}, 2^{n_i} - 1]$  and allows a similar range of representable numbers,  $2^{-16}$  to  $2^{16}$ , compared to the exponent bits in Float16. The fraction  $f$  is in  $[0, 1)$  and identically defined to the mantissa bits in floating-point numbers (without the hidden bit), which are encoded as the sum of powers of two with negative exponent. Multiplication, division, square root and power are easily implemented for a logarithmic number format with binary integer arithmetic and do not introduce any rounding errors unless the result is beyond the range of representable numbers. In contrast, addition and subtraction with logfixs is based on the Gaussian logarithms,

which often introduce a rounding error where floats avoid them due to their piecewise uniform distribution. For more details and a software implementation see `LogFixPoint16s.jl`<sup>41</sup>.

**Stochastic rounding.** The default rounding mode for floats, posits and logfixs is round-to-nearest<sup>10</sup>, which rounds an exact result  $x$  to the nearest representable number  $x_i$ . Stochastic rounding has recently emerged as an alternative rounding mode, beneficial for scientific computing<sup>31,33,42,43</sup>. In this rounding mode  $x$  is rounded down to a representable number  $x_1$  or up to  $x_2$  at probability proportional to the distance between  $x$  and  $x_1, x_2$

$$\text{round}_{\text{stoch}}(x) = \begin{cases} x_1 & \text{with probability } 1 - u^{-1}(x - x_1) \\ x_2 & \text{with probability } u^{-1}(x - x_1) \end{cases} \quad (3)$$

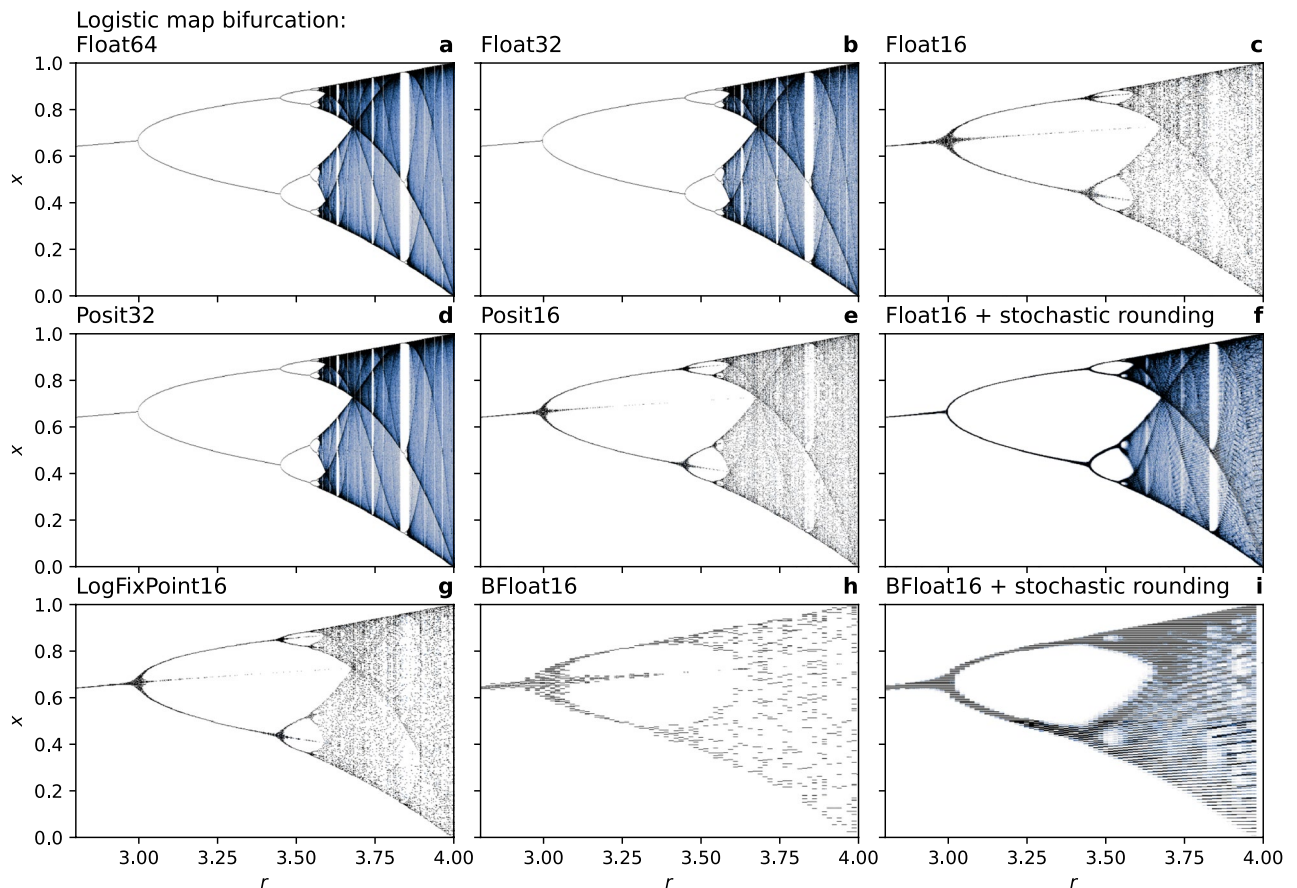
with  $u$  being the distance  $x_2 - x_1$  and  $x_1 \leq x \leq x_2$ . While deterministic rounding introduces a rounding error at most  $\pm \frac{u}{2}$ , the stochastic rounding error is bounded by  $\pm u$ . This occurs at low probability when  $x$  is very close to a representable number but rounded away from it. However, stochastic rounding is exact in expectation<sup>32,44</sup>, such that for repeated rounding the expectation of the rounded result converges to the exact  $x$ . Depending on the algorithm, the effective precision of a number format with stochastic rounding can therefore be higher than with deterministic round-to-nearest<sup>33,42</sup>, despite the same number of mantissa bits. For more details and the software implementation used here see the software package `StochasticRounding.jl`<sup>45</sup>.

### The logistic map with various number formats

The 1-variable logistic map is one of the most well-known chaotic maps, famous for its period-doubling bifurcation<sup>46</sup>. Starting with  $x_i \in [0, 1)$  at time step  $i = 0$  the logistic map with parameter  $r$  is

$$x_{i+1} = rx_i(1 - x_i) \quad (4)$$

For  $r = [1, 3]$  the system has a stable fixed point at  $\frac{r-1}{r}$  which all solutions will eventually converge to. For  $r = 3$  to  $r \approx 3.45$  the solution will oscillate between two values, for larger  $r$  this period doubles first to 4 time steps, then to 8, 16, ... and so on (Fig. 1a). For each period doubling the fixed points before the previous bifurcation



**Figure 1.** The bifurcation diagram of the logistic map simulated with various number formats, precisions and rounding modes. Bifurcation in the logistic map is a function of its parameter  $r$ . Shading represents a histogram of the solutions with darker colours denoting higher frequencies. Number formats, precisions and rounding modes in (a–i) as described in the respective titles.

continue to exist, but are now unstable. For  $r \approx 3.57$  and larger, chaotic solutions are found. However, intermittent ranges of  $r$  (often called islands of stability) exist which again have periodic solutions.

These properties of the logistic map are well simulated with Float64, Float32 and Posit32 (Fig. 1a,b,d) due to sufficiently high precision. However, at lower precision with Float16, Posit16, LogFixPoint16 and BFloat16 the unstable fixed points gain stability due to rounding errors (Fig. 1c,e,g,h, shown as third branch at bifurcations). Within the vicinity of these fixed points a diverging solution can be rounded back towards the unstable fixed points, which therefore become spuriously part of the simulated attractor. Furthermore, the chaotic solutions of the logistic map collapse into periodic orbits such that the dense attractor is only approximated by a low number of finite points (shown as sparsely dotted areas in the bifurcation diagram).

Stochastic rounding removes the aforementioned spurious stability of the unstable fixed points due to rounding errors at low precision (Fig. 1f,i). For Float16, stochastic rounding considerably improves the bifurcation diagram over deterministic rounding. Particularly the islands of stability reemerge and the chaotic solutions have a much denser attractor. While the bifurcation diagram with BFloat16 is poorly simulated with both rounding modes, stochastic rounding is still a clear improvement over deterministic rounding.

## Revisiting the generalised Bernoulli map

The generalised Bernoulli map<sup>11</sup> (also sometimes called the beta-shift<sup>47</sup> or the Renyi map<sup>48,49</sup>) is a 1-variable chaotic system starting with  $x_i \in [0, 1]$  at time step  $i = 0$  with the parameter  $\beta > 1$  defined as

$$x_{i+1} = f_\beta(x_i) = (\beta x_i) \bmod 1 \quad (5)$$

The modulo-operation mod satisfies that  $x \in [0, 1)$  at all future time steps. We note in passing that the generalised Bernoulli map is topologically conjugate to various other dynamical systems, including the aforementioned logistic map. Simulating this system with Float32 was found not to represent the periodic orbit spectrum well<sup>11</sup>, which is in turn closely related to the simulated invariant measure. For the generalised Bernoulli map the analytical invariant measure is known as<sup>50</sup>

$$h_\beta(x) = C \sum_{j=0}^{\infty} \beta^{-j\theta} \left( f_\beta^j(1) - x \right) \quad (6)$$

The Heaviside function is  $\theta$  and  $f_\beta^j(1)$  is the  $j$ -th time step of the Bernoulli map starting from  $x_0 = 1$ .  $C = 1$  is a normalisation constant, but for the calculation of Wasserstein distances (see section [Methods](#)) renormalization is applied so that the integral of  $h_\beta(x)$  over  $[0, 1]$  is equal to 1 and  $h_\beta(x)$  a probability density. To better visualise the invariant measure for varying  $\beta$  we introduce a normalisation  $\tilde{h}_\beta = h_\beta(x) / \max(h_\beta(x))$ , which is always in  $[0, 1]$  and can be applied to the analytical invariant measure as well as simulated ones.

We are revisiting the generalised Bernoulli map with various number formats and rounding modes to better understand a previously suggested pathology<sup>11</sup> as a function of arithmetic precision. While simulating the Bernoulli map numerically with a given number format, we perform both the multiplication and the subtraction in Eq. (5) with that format and avoid any conversion between number formats. This is in contrast to Boghosian et al. 2019, whose implementation converts  $x_i$  to Float64 before multiplication with  $\beta$  (as Float64) and possible subtraction with 1 (as Float64), i.e.  $x_{i+1} = \text{Float32}(\beta * \text{Float64}(x_i)) \bmod 1$ . While hardware allows for fused multiply-add operations without intermediate rounding error, similar to the conversion to Float64 here, the fused conditional subtraction in the modulo is generally not supported on hardware.

**The special  $\beta = 2$  case.** Boghosian et al. 2019 highlight that the Bernoulli map with  $\beta = 2$ , and similarly for every even integer, will collapse to  $x = 0$  after  $n$  time steps with any float format at arbitrary high but finite precision, where  $n$  is smaller than the number of bits. The subtraction with 1 acts as a bitshift towards more significant bits, pushing zero bits into the mantissa until  $x = 1$  and the modulo returns  $x = 0$  (Figure S3a, b and c). This phenomenon occurs as the Bernoulli map with  $\beta = 2$  (and similarly for larger even integers) does not introduce any arithmetic rounding error: Both the multiplication with  $\beta$  and subtraction with 1 are exact with floats (and also with posits). The multiplication with  $\beta = 2$  is exact as the base-2 exponent is simply increased by 1. The subtraction with 1 is exact as every finite positive float or posit can be written as  $2^e(1+f)$  for some integer  $e$  and a sum  $f \in [0, 1)$  of powers of two with negative exponents. Constraining the range to  $[1, 2)$ , where the subtraction is applied, yields  $e = 0$  and so subtracting 1 from the mantissa  $1+f$  is  $f$ , again a sum of powers of two, which is exactly representable with floats or posits.

The only occurring rounding error is in the initial conditions. While a randomly chosen  $x \in [0, 1)$  at infinite precision will have infinitely many non-zero mantissa bits, at finite precision those beyond the resolved mantissa bits are rounded to 0. Therefore, the least significant mantissa bit remains 0 after each iteration of the Bernoulli map while the same 0 bit from the previous iteration is further bit-shifted in. This behaviour holds for floats and posits, but it does not occur with logfixs. All multiplications are exact with logfixs (unless under or overflows occur), but in contrast to floats and posits a rounding error occurs in the subtraction, with the possibility of setting the least significant mantissa bit to 1. This rounding error is effective at preventing a collapse of the attractor (Figure S3d). The Bernoulli map with  $\beta = 2$  and simulated with floats or posits is therefore special, as it is a chaotic system that does not involve any arithmetic rounding errors beyond the rounding of the initial conditions. However, the simulation of most other systems, including the generalised Bernoulli map with  $1 < \beta < 2$ , involves rounding errors with any finite precision number format.

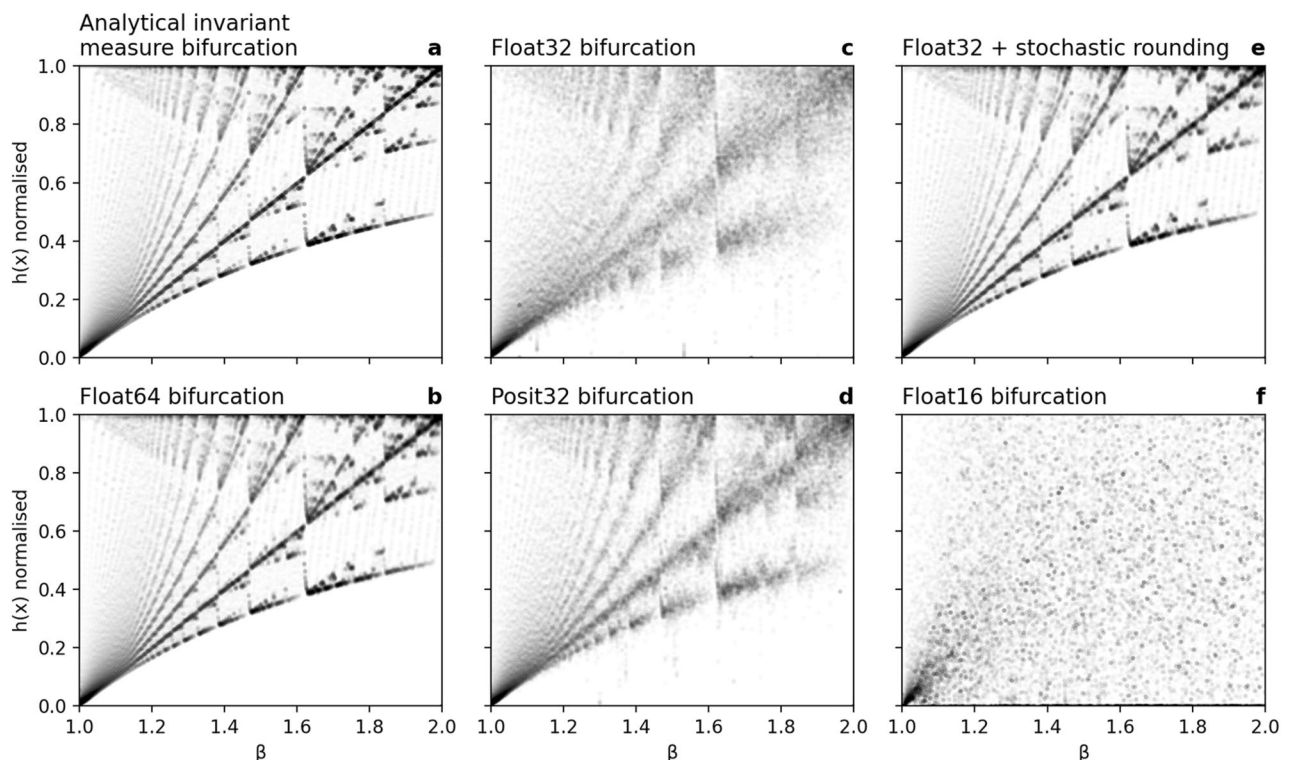


**Bifurcation of the invariant measure.** Given that the analytical invariant measure is known for the generalised Bernoulli map (Eq. 6), we can assess its representation with various number formats at different levels of precision. Boghosian et al. 2019 conclude that the invariant measure with Float32 is an inaccurate approximation of the analytical invariant measure. While this difference is even more pronounced with Float16 arithmetic (Figure S2), with Float64 the invariant measure is comparably accurate. The question therefore arises as to whether the discrepancy of the invariant measures vanishes with higher precision, or whether a pathology persists at any precision level for some  $\beta$ .

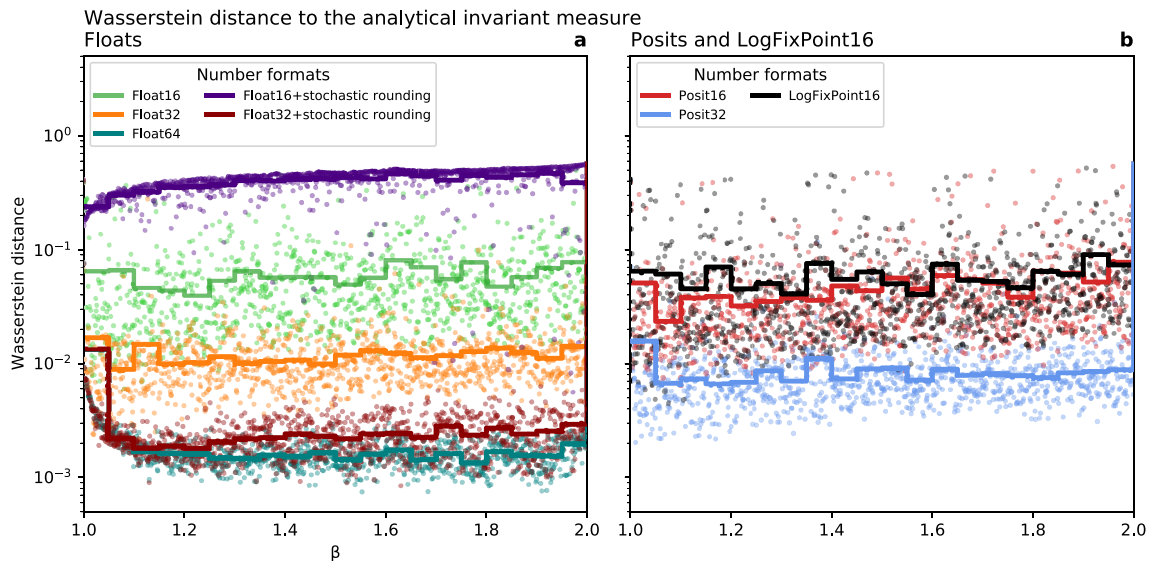
The analytical invariant measure of the generalised Bernoulli map consists of many step functions taking values on a discrete set of points (Figure S2), which bifurcate with increasing  $\beta$  (Fig. 2a). These “quantization levels” cannot be exactly represented with Float16 or Float32 arithmetic (Figure S2), such that their bifurcation is visually blurred (Fig. 2c). The visually sharp representation of this bifurcation of the quantization levels with Float64 indicates a much more accurate approximation to the analytical invariant measure (Fig. 2b). Due to the higher precision of 32-bit posit arithmetic (Posit32) around  $\pm 1$ , the bifurcation is slightly improved with Posit32 over Float32 (Fig. 2d). Given the inaccurate representation of the invariant measure with Float16 (Figure S2), its bifurcation has little resemblance to the analytical bifurcation (Fig. 2f).

**Effects of stochastic rounding.** Augmenting Float32 with stochastic rounding considerably improves the bifurcation (Fig. 2e) and makes it virtually indistinguishable from Float64 or the analytic bifurcation. However, stochastic rounding does not decrease the rounding error accumulated over many iterations in a forecast (Figure S4), such that the effective precision is not increased over deterministic rounding. But introducing stochasticity prevents the convergence onto periodic orbits, which are otherwise present with deterministic rounding<sup>11</sup>. Previously inaccessible regions of the attractor can be reached with stochastic rounding as the simulation is frequently pushed off any periodic orbit. This advantage of stochastic rounding is also observed in the logistic map (section “The logistic map with various number formats”). While periodic orbits are not fully removed from solutions due to the use of pseudo random number generators (PRNG) that are themselves periodic, the periods of PRNGs are usually so long that effectively any periodicity is avoided. The period length of Mersenne Twister<sup>51</sup>, the most widely-used PRNG, is  $2^{19937} - 1$  and still sufficiently long with  $2^{128} - 1$  for the faster Xoroshiro128 + <sup>52,53</sup> PRNG that is used in StochasticRounding.jl.

The agreement of the analytical and simulated invariant measures is quantified with the Wasserstein distance (section Methods). For  $\beta = 2$  the analytical invariant measure is the uniform distribution  $U(0, 1)$ , whereas all float and posit formats for both deterministic and stochastic rounding simulate a collapse of the attractor to zero such that the invariant measure is the Dirac delta distribution (Fig. 3). The Wasserstein distance  $W_1$  is in all these cases  $W_1 = 0.5$  and does not improve with precision. However, as previously mentioned, the rounding errors



**Figure 2.** Bifurcation of the quantization levels corresponding to the invariant measure in the generalised Bernoulli map as simulated with various number formats. (a) Analytical bifurcation  $h_\beta(x)$  from the exact invariant measure, normalised by  $\max(h_\beta(x))$ , compared to the invariant measure by simulating the Bernoulli map with (b) Float64, (c) Float32, (d) Posit32, (e) Float32 + stochastic rounding, and (f) Float16.



**Figure 3.** Agreement between the simulated and analytical invariant measures in the generalised Bernoulli map quantified by the Wasserstein distance. For all values  $1 \leq \beta < 2$  a higher precision number format yields a better agreement with the analytical Bernoulli map. Simulations using (a) Floats with and without stochastic rounding, (b) Posits and logarithmic fixed-point numbers. The Wasserstein distances are calculated for the invariant measures obtained from  $N = 10^3$  simulations for each value of  $\beta$ . Scatter points denote individual Wasserstein distances, solid lines indicate averages across a range of  $\beta$  as indicated by steps.

from logfixs prevent a collapse such that for LogFixPoint16 the invariant measure is much better approximated, with  $W_1 = 0.05$ .

For  $1 \leq \beta < 2$  the Wasserstein distance is always reduced going to higher precision, supporting the inference that only the case  $\beta = 2$  (and other even integers) presents a pathology where higher precision does not improve the simulated invariant measure arising from the generalised Bernoulli map. The Wasserstein distances of Float32 with stochastic rounding are similarly low to Float64 and no significant difference can be found. However, using stochastic rounding with Float16 is worse than deterministic rounding as here the stochasticity makes it possible that the invariant measure of the generalised Bernoulli map collapses to 0, which is a fixed point (Figure S5). For Float32 the probability of such an occurrence is low and is not observed here. With Float16 most simulations collapse within a few thousand iterations transforming their invariant measures into Dirac distributions. Whether this problem generalises to other systems is questionable. We suspect that this may be a feature of low-dimensional dynamics and may not arise commonly in higher dimensional systems: There the chance of a stochastic perturbation onto a fixed point becomes vanishingly small even at very low precision. Other natural systems do not have fixed points due to time-dependent forcing. The posit format is slightly better than floats at both 16 and 32-bit, as expected from the slightly higher precision.

### Orbits in the Lorenz 1996 system

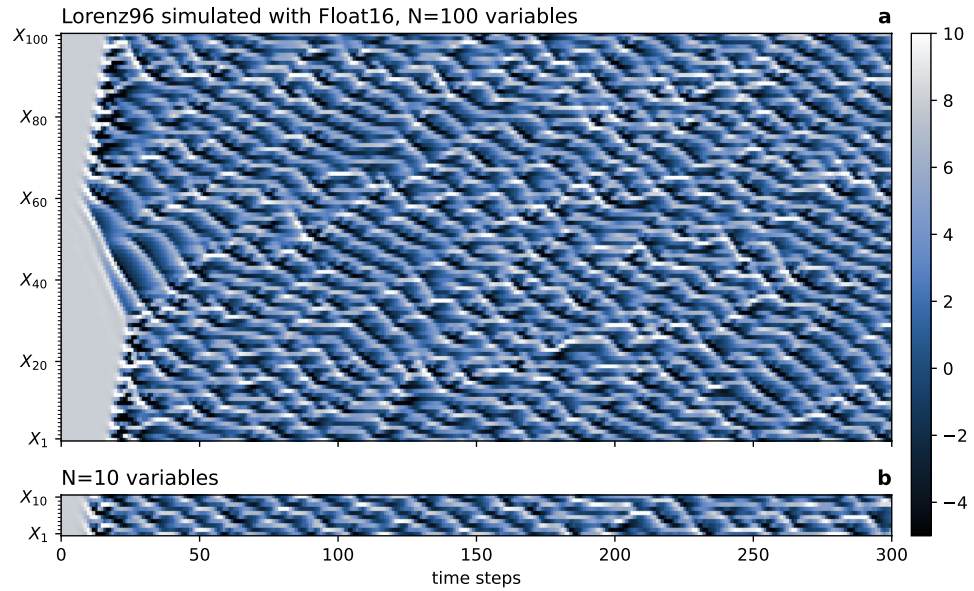
In contrast to the 1-variable generalised Bernoulli map, most continuous natural systems are simulated numerically with as many variables as computationally affordable by increasing resolution and/or complexity. Weather forecast and climate models often use millions of independent variables that result from a discretisation of continuous variables on the globe. While short periodic orbits in low precision are problematic in the simulation of few-variable systems as discussed in the previous sections, this section tests the hypothesis that large systems are unaffected for all practical purposes.

To investigate the dependence of periodic orbits on the number of variables in the system we consider the chaotic Lorenz 1996 system<sup>54,55</sup>. This system has been widely studied for data assimilation<sup>56</sup> and machine learning<sup>57,58</sup>. With  $N$  variables  $X_i, i = 1, \dots, N$  the one-layer version is a system of coupled ordinary differential equations.

$$\frac{dX_i}{dt} = X_{i-1}(X_{i+1} - X_{i-2}) - X_i + F \quad (7)$$

in a one-dimensional spatial domain with periodic boundary conditions,  $X_{N+1} = X_1$  etc. The term  $X_{i-1}(X_{i+1} - X_{i-2})$  implements nonlinear advection, and drag is represented with the relaxation term  $-X_i$ . The forcing  $F$  is the single parameter in the Lorenz 1996 system fixed at the common default  $F = 8$  which produces chaotic solutions. The forcing is steady in time and constant in space. The system exhibits dynamics of nonlinear wave-wave interactions (Fig. 4a), which are reasonably independent of the number of variables (Fig. 4b). The system can be integrated with as little as  $N = 4$  variables without an obvious degradation of the simulated dynamics.

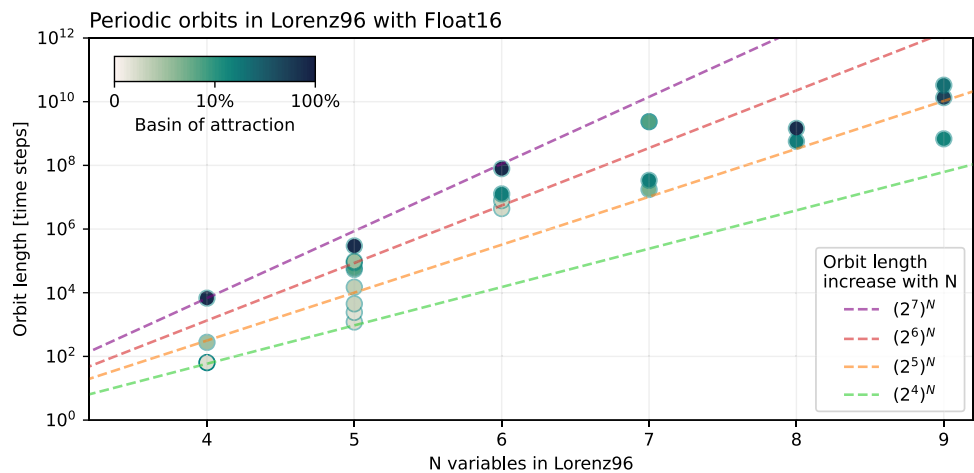
The initial conditions are in equilibrium  $X_i = F, i \neq j \forall i$  with only a single variable which is slightly perturbed  $X_j = F + 0.005 + 0.01\varepsilon$  with  $\varepsilon \sim U(0, 1)$ , drawn from a random uniform distribution in  $[0, 1)$ . Due to the periodic boundary conditions and the spatially constant forcing the system is spatially invariant. After



**Figure 4.** The Lorenz 1996 system simulated with Float16 arithmetic. A Hovmoeller diagram visualising the temporal evolution of every variable  $X_i$  encoded in colour. (a)  $N = 100$  variables starting from equilibrium  $X_i = 8$  (coloured as grey) with a small perturbation in  $X_{50}$ , and (b) same as (a) but  $N = 10$  variables, starting with a perturbation in  $X_5$ . Pixels visible in the shading represent individual variables and time steps.

some several hundred time steps, the information from the initial conditions is removed through the chaotic dynamical evolution (Fig. 4a). The invariant measure  $\mu$  of one variable  $X_i$  is therefore identical to that of any other  $\mu(X_i) = \mu(X_j) \forall i, j$ , as will be further discussed below.

The Lorenz 1996 system is discretized in time using the 4-th order Runge–Kutta scheme<sup>59</sup> with a time step of  $\Delta t = 0.1$ . At this temporal resolution the system can also be integrated using a low-precision number format such as Float16 (Fig. 4). For more details and a software implementation see Lorenz96.jl<sup>60</sup>. Time integration scheme and time step also have a large impact on chaotic trajectories and hence on the periodic orbit analysis presented here. Numerical stability and performance usually dictate this choice, but low precision can add an additional constraint: Using a shorter time step can cause stagnation as tendencies are too small to be added in the time integration. Stagnation from low precision can be overcome with a compensated time integration<sup>61</sup> or with stochastic rounding<sup>32</sup>.



**Figure 5.** Periodic orbits in Lorenz96 simulated with Float16 with an increasing number of variables. Initial conditions are randomly taken from a high-precision simulation. Basins of attraction (shading) correspond to the share of initial conditions that converge to the respective periodic orbit. Dashed lines provide an orientation for the exponential increase in orbit lengths with the number of variables.



**Longer orbits with more variables.** Using  $N = 4$  variables in the Lorenz 1996 system simulated with Float16, the longest periodic orbit we find is 6756 time steps long (Fig. 5 and Table S1). The basin of attraction is about 0.82, meaning that about 82% of the randomly chosen initial conditions converge onto this orbit. Increasing the number of variables to  $N = 5$ , the longest periodic orbit found increased to a length of 294,995 time steps at a similarly large basin of attraction. For  $N > 9$  the orbit search becomes computationally very demanding and requires more than several days on sizable compute clusters with 100 cores (see section **Methods** for a description of how the orbit search is parallelised). For  $N = 9$  though, the longest periodic orbit we were able to find has a period of 32,930,252,532 time steps. For a list of periodic orbits found in the Lorenz 1996 system and their respective minimums see Table S1.

**More variables instead of higher precision.** In most cases between 4 and 9 variables in the Lorenz 1996 system, the longest orbit is also the one with the largest basin of attraction. The longer the orbit the larger the occupied state space of possible values the variables  $X_i$  can take at a given precision. Consequently, the assumption is that it is most likely that a given trajectory ends up on the longest orbit. However, we also found a counter example as the longest orbit with  $N = 8$  variables is shorter than the longest with 9 variables (Fig. 5).

The orbit length increases approximately exponentially following a scaling of about  $16^N$  to  $128^N$  from  $N = 4$  to  $N = 9$ . Such an exponential increase translates to about 4 to 7 effective bits of freedom (as  $2^4 = 16$ ,  $2^7 = 128$ ) for every additional variable in Lorenz 1996 represented with Float16. However, the computational resources limit the orbit search for larger  $N$ , making it hard to constrain this exponential scaling further. Assuming a similar exponential orbit increase holds for larger  $N$ , extrapolation of these findings suggests orbit lengths on the order of about  $10^{100,000}$  for million-variable systems simulated with Float16. This is far beyond the reach of any computational resources currently available. In that sense, while a simulation of such large systems would eventually be periodic, a periodic solution will never be reached.

Longer orbits are promising to avoid periodic solutions in low precision, but short periodic orbits do not necessarily misrepresent a reference invariant measure. We assess the agreement of invariant measures using the Wasserstein distance as before. As a reference invariant measure  $\mu(X_{ref})$  we integrate the Lorenz 1996 system for 1,000,000 time steps with  $N = 500$  variables using Float64 arithmetic. The Wasserstein distance is then  $W(\mu(X_{ref}), \mu(X))$  with  $X$  representing the variables from a Lorenz 1996 simulation with  $N$  variables using either Float16 or Float64 arithmetic.

Using only  $N = 4$  variables in the simulation of Lorenz 1996 yields an invariant measure with little resemblance to the reference (Fig. 6a,f), regardless of the number format. While more variables yield an invariant measure that converges to the reference, there is virtually no difference whether Float16 or Float64 arithmetic is used (Fig. 6b–e). The Wasserstein distance significantly reduces with an increasing number of variables, but not with higher precision (Fig. 6g). Given a certain availability of computational resources a better invariant measure is therefore obtained by reducing the precision and reinvesting the performance gain into more variables.

**Discussion.** From a rigorous mathematical perspective, we would like to know how the periods of orbits increase with the number of variables and precision. Experimentally, we find orbits that exponentially increase in length with the number of independent variables in the Lorenz 1996 system. Every variable is found to add between 4 and 7 maximum entropy bits that extend orbits by a factor of  $2^4$  to  $2^7$ . Similarly, assuming maximum entropy for additional mantissa bits, the expected length of periodic orbits doubles with every additional mantissa bit in precision.

The Lorenz 1996 system is deemed to be *too easy*<sup>62</sup> as a toy model for machine learning, meaning that it is too homogeneously chaotic to be actually a challenging problem. In our case, the conclusions from Lorenz 1996 may not translate directly to more chaotic systems: the exponential growth of the periodic orbit length with the number of variables in such systems might be weaker.

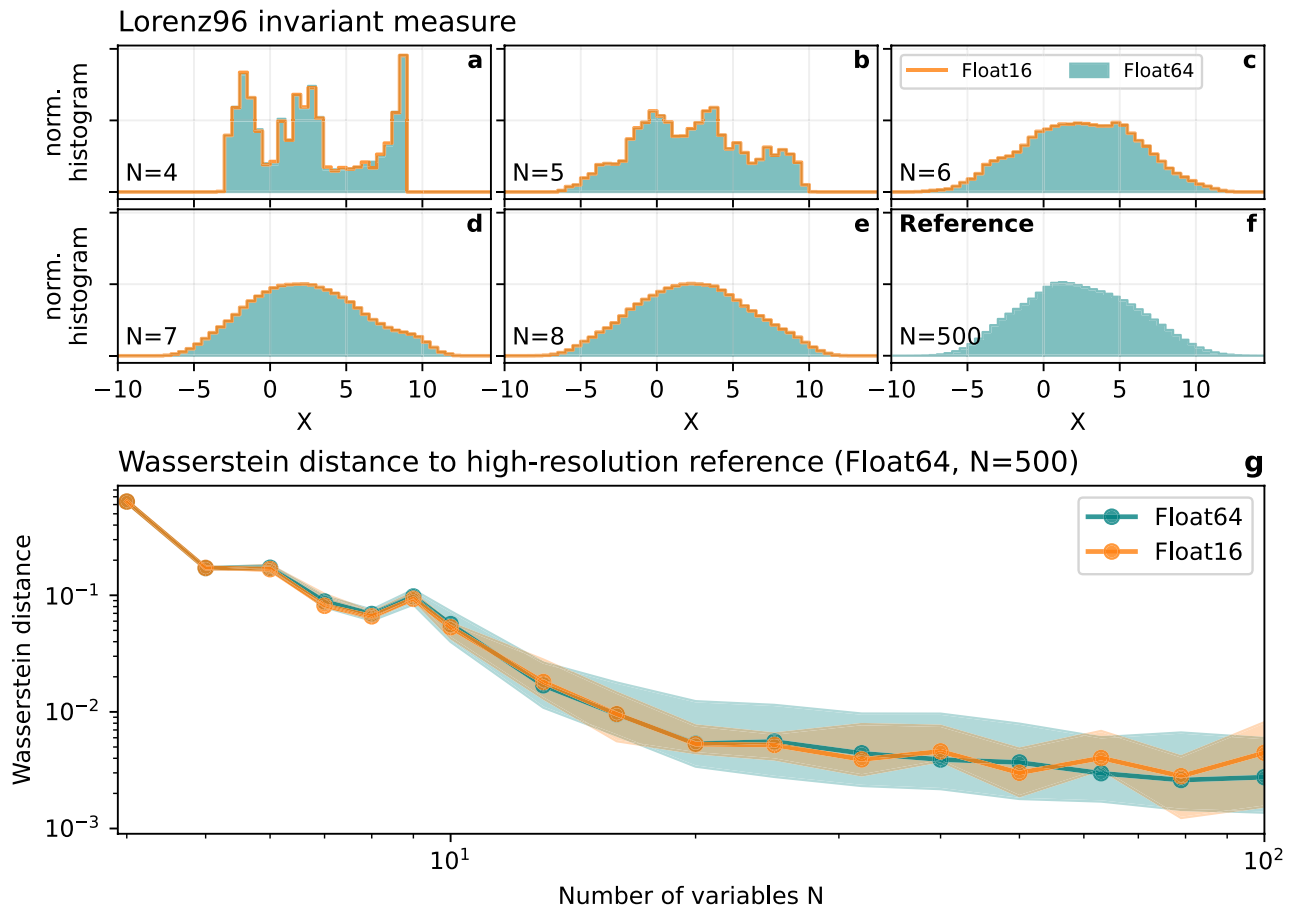
For more complex simulations of natural systems such as million-variable climate models, the periodic orbits found here would naively extrapolate to lengths beyond  $10^{100,000}$  time steps even in 16-bit precision. The largest orbits will very likely remain out of reach on future generations of supercomputers, especially if the performance gained from low-precision simulations is reinvested into higher resolution or complexity. In the context of climate models, this supports a vision of low-precision but high-resolution simulations with added stochasticity to accelerate and improve climate predictions.

## Conclusions

We analysed the bifurcation in the logistic map with different number formats and rounding modes. The rounding errors from 16-bit floats, posits or logfixs with deterministic rounding can spuriously stabilise the unstable fixed points. Furthermore, the dense attractors collapse into short periodic orbits. Both issues can be addressed with stochastic rounding. It prevents periodic orbits as a chaotic trajectory is regularly pushed off any periodic orbit due to the stochastic perturbation in rounding. This considerably improves the logistic map bifurcation with Float16.

The periodic orbit spectrum in the generalised Bernoulli map was analysed with different number formats and levels of arithmetic precision. While there are very special cases (such as  $\beta = 2$ ) in which the system's simulation is greatly degenerated at any precision, in all other cases simulations were found to improve with higher precision. 16 and 32-bit arithmetic result in short periodic orbits in the Bernoulli map, but a sufficiently high precision reduces the error in the invariant measure to a minimum.

Stochastic rounding is also found to be especially beneficial for 32-bit floats in the Bernoulli map. Simulated invariant measures are improved as chaotic trajectories travel with stochastic rounding through otherwise unreachable state space. However, in the Bernoulli map stochastic rounding also causes a non-zero chance that



**Figure 6.** Improvement of the simulated Lorenz 1996 invariant measure with increasing number of variables. (a) The invariant measure of Lorenz 1996 simulated with  $N = 4$  variables. Using Float16 (orange line outlining the histogram) or Float64 (teal filled histogram) arithmetic yields a virtually identical invariant measure. (b–e) as (a) but with an increasing number of variables. (f) The reference invariant measure obtained from  $N = 500$  variables using Float64 arithmetic. (g) The Wasserstein distance of the simulated Lorenz 1996 system with respect to the reference. Invariant measures are taken from all available variables, which are invariant due to periodic boundary conditions and spatially-independent forcing (see Eq. 7). Shadings in (g) represent the 5–95% confidence interval and solid lines the median obtained from an ensemble simulation with 100 members starting from slightly different random initial conditions.

the system collapses on the fixed point. But for any precision higher than Float16 and in systems with more variables this chance quickly vanishes. In many complex natural systems the fixed points are not close to the attractor, further limiting the relevance of this issue in practice.

Increasing the number of variables in the Lorenz 1996 system, we find that more variables improve the simulated invariant measure much more than increased precision with fewer variables does: Doubling the amount of variables yields a more accurate invariant measure than doubling the precision of a high-resolution and high-precision reference. This provides evidence that computational resources should be invested in higher resolution rather than higher precision for the simulation of continuous chaotic systems.

## Methods

**An improved random number generator for uniformly distributed floats.** Conventional random number generation for a float  $f$  from a uniform distribution  $U(0, 1)$  uses the following technique: First, 10/23/52 random bits (for Float16/Float32/Float64 respectively) from an unsigned integer are used to set the mantissa bits of floating-point 1. This creates a floating-point number that is uniformly distributed in  $[1, 2)$  as all float formats are uniformly distributed in that range (the exponent bits are constant). Second, 1 is subtracted to obtain a float in  $[0, 1)$ , i.e.  $f \sim U(1, 2) - 1$ . While this approach is fast, it is statistically imperfect as the resulting distribution  $U(1, 2) - 1$  does not contain all floats in  $[0, 1)$ . There are  $2^{23}$  Float32s in  $[1, 2)$  but the subtraction maps those only to a subset of all  $1,065,353,216 \approx 2^{30}$  Float32s in  $[0, 1)$ . This technique only samples from every second float in  $[\frac{1}{2}, 1)$ , every fourth in  $[\frac{1}{4}, \frac{1}{2})$  and so every  $2n$ -th float in  $[2^{-n}, 2^{-n+1})$ . Furthermore, the smallest positive number that can be obtained is about  $10^{-7}$  for Float32 and  $10^{-16}$  for Float64. This is many orders of magnitude larger than minpos, the smallest representable positive float, which is about  $10^{-45}$ ,  $10^{-324}$  for Float32, Float64, respectively.

We therefore developed a statistically improved conversion from a random unsigned integer to a uniformly distributed float in  $[0, 1)$ . Counting the number of leading zeros  $l$  of a random unsigned integer yields  $l = 0$  at probability  $\frac{1}{2}$ ,  $l = 1$  at probability  $\frac{1}{4}$  and  $l = k$  at probability  $2^{-k-1}$ . These probabilities correspond exactly to the share of power-2 exponents in the unit range  $[0, 1)$ . Consequently, we translate the number of leading zeros  $l$  to the respective exponent bits and use the remaining bits of the unsigned integers for the mantissa bits. The statistical flaws from the conventional conversion as presented above are avoided, but for practical reasons the smallest float that can be sampled is about  $10^{-20}$  for both Float32 and Float64. It is therefore practically impossible to sample a zero with this technique, which is in contrast to the conventional technique. Just as the chance of obtaining a zero in  $[0, 1)$  is 0 for real numbers, the method here also has a zero chance for floats. We implement this method in the software package `RandomNumbers.jl` and use it throughout this study.

**Monte Carlo orbit search.** For the generalised Bernoulli map, the random number generator described above is used to sample from all floats in  $[0, 1)$  to obtain a representative subset of all initial conditions. While there is no guarantee that all orbits are found, those found have the largest basin of attraction. While it is easily possible to miss a periodic orbit, those missed are expected to have a very small basin of attraction and therefore a negligible contribution to the invariant measure. Estimating the invariant measure from these orbits is therefore also expected to be an unbiased approximation that converges to the exact invariant measure. The exact invariant measure is obtained by finding all simulated orbits and their exact basins of attraction rather than using a random set of initial conditions. We verify this methodology for Float16 and Float32, where the exact invariant measure can be calculated in Figure S1. While we cannot find all orbits with Float64, the Monte Carlo-based invariant measure converges to the analytical invariant measure and is for the same sample size a better approximation than using Float16 or Float32. Despite the high precision, a Float64 simulation of the generalised Bernoulli map still substantially degrades some properties of the analytical system: The topological entropy, measuring how trajectories diverge onto distinct orbits, is positive in the analytical system, representing chaotic solutions. However, even with the high precision of Float64 the topological entropy is negative, as trajectories eventually converge onto periodic orbits.

For the Lorenz 1996 system, the space of all possible initial conditions is much larger than the space the attractor occupies. It is therefore more efficient to only choose initial conditions randomly that are already part of, or at least close to, the attractor. The basin of attraction here means therefore the relative share of the initial conditions from the attractor that end up on a given orbit, and not from all possible initial conditions. To obtain an initial condition for the Lorenz 1996 system, one first starts a high-precision simulation from a given initial condition including a small stochastic perturbation (see section ["Orbits in the Lorenz 1996 system"](#) for more details). After disregarding a spin-up the information about the chosen initial condition is removed and the stochastic perturbation grows into a fully independent random initial condition. Converting a random time step from a high-precision simulation into the given number format is then intended to emulate sampling from the invariant measure at low precision.

**Efficient orbit search with distributed computing.** To find an orbit in a simulation, Eq. (1) is used after every time step to check for equality with a previous time step. However, before an orbit is found, it is unknown whether a given initial condition  $X_{t_0}$  is already part of the orbit or still part of the trajectory that is yet to converge onto an orbit. It is possible to use the last time step of a very long spin-up simulation as  $X_{t_0}$ . This strategy limits the chance that  $X_{t_0}$  is not yet part of the orbit, but does not provide a guarantee, nor is it efficient. Instead we implemented a strategy whereby  $X_{t_0}$  is updated during simulation and slowly moves forward in time: Updates like  $t_0 = \text{round}(\sqrt{t_1})$  or  $t_0 = \text{round}(\log(t_1))$ , with  $t_0 < t_1$  integers, indicating the time steps, are used. In particular, we use several past time steps of the simulation as  $X_{t_0}$  to check for periodicity. Checking for periodicity with *all* past time steps is inefficient as they would have to be stored and  $O(t_1^2)$  checks have to be performed in total for all time steps from  $t_0$  to  $t_1$ . In contrast, for a constant number of checks per time step, the total number of checks increases only linearly with the simulation time.

Finding  $n$  orbits from  $n$  different initial conditions is a problem that is parallelizable into  $n$  independent processes calculated on  $n_p \leq n$  processors. We follow ideas of the MapReduce framework: Each worker process starts with a different initial condition and simulates the dynamical system independently of other processes until an orbit is found. This orbit is passed to the main process, which reduces successively all  $n$  orbits found into a list of unique orbits, as several initial conditions can yield the same orbit. Instead of defining an orbit by all of the points on it, which would be computationally inefficient for very long orbits, we describe an orbit by the period length and its minimum. The minimum of an orbit is the point for which the  $L^2$  norm is minimised. While it is theoretically possible that an orbit has several minima with identical norms, this occurred rarely in our applications. Two such orbits that are falsely identified as respectively unique are then merged in post-processing. A uniqueness check between two orbits (or one orbit and a list of orbits, in which case the uniqueness check is pairwise against every orbit in the list) is unsuccessful and yields a single orbit only if all of the three following criteria are fulfilled: 1) Length: the two orbits must have the same period length; 2) Minimum norm: the norms of the orbits' minima have to be identical; 3) Minimum: the orbits' minima, including possible rotation of the variables for spatially periodic solutions, have to be bitwise identical. While criterion 3 is sufficient to identify the uniqueness of two orbits, it is computationally more efficient to check for criterion 3 only if criterion 2 is fulfilled, which is only checked if criterion 1 is fulfilled, hence the proposed order.

**Wasserstein distance.** The invariant measure of a chaotic dynamical system is estimated with histogram binning. To assess the agreement of two histograms representing invariant measures (either simulated or analytical) we use the Wasserstein distance, a metric that derives from the theory of optimal transport, with an  $L^1$

cost. The Wasserstein distance is defined as the least cost at which one can transport all probability mass from histogram  $\mu$  to another histogram  $\nu$ , where the cost to move mass  $m$  from a bin of  $\mu$  at location  $x$  to a bin at location  $y$  of  $\nu$  is  $m|x - y|^{3,63}$ . This gives a non-parametric method to compare probability distributions which accounts for both differences in the probabilities of events as well as their separations in the underlying space, so that closeness in Wasserstein distance truly corresponds to a natural notion of closeness between probability distributions<sup>63</sup>, Thm 7.12.

## Data availability

The repository BernoulliMap is available at <https://github.com/milankl/BernoulliMap><sup>64</sup> and contains the software to produce the analysis presented here. Lorenz96.jl (v0.3) is available at <https://github.com/milankl/Lorenz96.jl><sup>60</sup>. LogFixPoint16s.jl (v0.3) is available at <https://github.com/milankl/LogFixPoint16s.jl><sup>41</sup>. StochasticRounding.jl (v0.6) is available at <https://github.com/milankl/StochasticRounding.jl><sup>45</sup>.

Received: 31 October 2022; Accepted: 14 June 2023

Published online: 14 July 2023

## References

- Bauer, P., Thorpe, A. & Brunet, G. The quiet revolution of numerical weather prediction. *Nature* **525**, 47–55. <https://doi.org/10.1038/nature14956> (2015).
- Palmer, T. The ECMWF ensemble prediction system: Looking back (more than) 25 years and projecting forward 25 years. *Q. J. R. Meteorol. Soc.* **145**, 12–24 (2019).
- Cummings, R. M., Mason, W. H., Morton, S. A. & McDaniel, D. R. *Applied Computational Aerodynamics: A Modern Engineering Approach*. (Cambridge University Press, 2015).
- Moran, J. *An Introduction to Theoretical and Computational Aerodynamics*. (Courier Corporation, 2003).
- Cornish, N. J. Chaos and gravitational waves. *Phys. Rev. D* **64**, 084011 (2001).
- Springel, V. The cosmological simulation code gadget-2. *Mon. Not. R. Astron. Soc.* **364**, 1105–1134 (2005).
- Coveney, P. V. & Wan, S. On the calculation of equilibrium thermodynamic properties from molecular dynamics. *Phys. Chem. Chem. Phys.* **18**, 30236–30240 (2016).
- Mazzi, S. *et al.* Enhanced performance in fusion plasmas through turbulence suppression by megaelectronvolt ions. *Nat. Phys.* **18**, 776–782 (2022).
- Ricci, P. *et al.* Simulation of plasma turbulence in scrape-off layer conditions: the GBS code, simulation results and code validation. *Plasma Phys. Control. Fusion* **54**, 124047 (2012).
- IEEE. IEEE Standard for Binary Floating-Point Arithmetic. *ANSI IEEE Std 754–1985 1–20*. <https://doi.org/10.1109/IEEESTD.1985.82928> (1985).
- Boghossian, B. M., Coveney, P. V. & Wang, H. A new pathology in the simulation of chaotic dynamical systems on digital computers. *Adv. Theory Simul.* **2**, 1900125 (2019).
- Cvitanović, P. Periodic orbits as the skeleton of classical and quantum chaos. *Phys. Nonlinear Phenom.* **51**, 138–151 (1991).
- Lasagna, D. Sensitivity of long periodic orbits of chaotic systems. *Phys. Rev. E* **102**, 052220 (2020).
- Leboeuf, P. Periodic orbit spectrum in terms of Ruelle–Pollicott resonances. *Phys. Rev. E Stat. Nonlinear Soft Matter Phys.* **69**, 026204 (2004).
- Ruelle, D. & Takens, F. On the nature of turbulence. *Rencontres Phys.-Mathématiciens Strasbg.-RCP25* **12**, 1–44 (1971).
- Maiocchi, C. C., Lucarini, V. & Gritsun, A. Decomposing the dynamics of the Lorenz 1963 model using unstable periodic orbits: Averages, transitions, and quasi-invariant sets. *Chaos Interdiscip. J. Nonlinear Sci.* **32**, 033129 (2022).
- Eckmann, J.-P. & Ruelle, D. Ergodic theory of chaos and strange attractors. in *The Theory of Chaotic Attractors* (eds. Hunt, B. R., Li, T.-Y., Kennedy, J. A. & Nusse, H. E.) 273–312 (Springer, 2004). [https://doi.org/10.1007/978-0-387-21830-4\\_17](https://doi.org/10.1007/978-0-387-21830-4_17).
- Cvitanović, P. Recurrent flows: The clockwork behind turbulence. *J. Fluid Mech.* **726**, 1–4 (2013).
- Ghil, M. & Malanotte-Rizzoli, P. Data Assimilation in Meteorology and Oceanography. in *Advances in Geophysics* (eds. Dmowska, R. & Saltzman, B.) vol. 33 141–266 (Elsevier, 1991).
- Butcher, J. C. *Numerical Methods for Ordinary Differential Equations*. (Wiley, 2016).
- Higham, N. J. *Accuracy and stability of numerical algorithms*. (SIAM, 2002).
- Fuhrer, O. *et al.* Near-global climate simulation at 1 km resolution: Establishing a performance baseline on 4888 GPUs with COSMO 5.0. *Geosci. Model Dev.* **11**, 1665–1681 (2018).
- Nakano, M., Yashiro, H., Kodama, C. & Tomita, H. Single precision in the dynamical core of a nonhydrostatic global atmospheric model: Evaluation using a baroclinic wave test case. *Mon. Weather Rev.* **146**, 409–416 (2018).
- Vaña, F. *et al.* Single precision in weather forecasting models: An evaluation with the IFS. *Mon. Weather Rev.* **145**, 495–502 (2017).
- Markidis, S., Chien, S. W. D., Laure, E., Peng, I. B. & Vetter, J. S. NVIDIA Tensor Core Programmability, Performance Precision. in *2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)* 522–531 (2018). <https://doi.org/10.1109/IPDPSW.2018.00091>.
- Jouppi, N., Young, C., Patil, N. & Patterson, D. Motivation for and evaluation of the first tensor processing unit. *IEEE Micro* **38**, 10–19 (2018).
- Odajima, T. *et al.* Preliminary performance evaluation of the Fujitsu A64FX Using HPC Applications. in *2020 IEEE International Conference on Cluster Computing (CLUSTER)* 523–530 (2020). <https://doi.org/10.1109/CLUSTER49012.2020.00075>.
- Sato, M. *et al.* Co-design for A64FX manycore processor and “Fugaku”, in *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis* 1–15 (2020). <https://doi.org/10.1109/SC41405.2020.00051>.
- Gustafson, J. & Yonemoto, I. Beating Floating point at its own game: Posit arithmetic. *Supercomput. Front. Innov.* **4**, 16 (2017).
- Johnson, J. Efficient, arbitrarily high precision hardware logarithmic arithmetic for linear algebra. in *2020 IEEE 27th Symposium on Computer Arithmetic (ARITH)* 25–32 (2020). <https://doi.org/10.1109/ARITH48897.2020.00013>.
- Hopkins, M., Mikaitis, M., Lester, D. R. & Furber, S. Stochastic rounding and reduced-precision fixed-point arithmetic for solving neural ordinary differential equations. *Philos. Trans. R. Soc. Math. Phys. Eng. Sci.* **378**, 20190052 (2020).
- Croci, M., Fasi, M., Higham, N. J., Mary, T. & Mikaitis, M. Stochastic rounding: Implementation, error analysis and applications. *R. Soc. Open Sci.* **9**, 211631.
- Paxton, E. A., Chantry, M., Klöwer, M., Saffin, L. & Palmer, T. Climate modeling in low precision: Effects of both deterministic and stochastic rounding. *J. Clim.* **35**, 1215–1229 (2022).
- Urmitsky, D. J. Shadowing unstable orbits of the Sitnikov elliptic three-body problem. *Mon. Not. R. Astron. Soc.* **407**, 804–811 (2010).
- Yahniz, G., Hof, B. & Budanur, N. B. Coarse graining the state space of a turbulent flow using periodic orbits. *Phys. Rev. Lett.* **126**, 244502 (2021).



36. Viswanath, D. Recurrent motions within plane Couette turbulence. *J. Fluid Mech.* **580**, 339–358 (2007).
37. IEEE. IEEE Standard for Floating-Point Arithmetic. *IEEE Std 754–2008* 1–70 (2008). <https://doi.org/10.1109/IEEESTD.2008.4610935>.
38. Gustafson, J. L. *The End of Error: Unum Computing*. (Chapman and Hall/CRC, 2015).
39. Klöwer, M., Düben, P. D. & Palmer, T. N. Posits as an alternative to floats for weather and climate models. in *Proceedings of the Conference for Next Generation Arithmetic 2019 on CoNGA'19* 1–8 (ACM Press, 2019). <https://doi.org/10.1145/3316279.3316281>.
40. Klöwer, M., Düben, P. D. & Palmer, T. N. Number formats, error mitigation, and scope for 16-bit arithmetics in weather and climate modeling analyzed with a shallow water model. *J. Adv. Model. Earth Syst.* **12**, e2020MS002246 (2020).
41. Klöwer, M. LogFixPoint16s.jl: A 16-bit logarithmic fixed-point number format (v0.3). *Zenodo*. <https://doi.org/10.5281/zenodo.8138366> (2023).
42. Croci, M. & Giles, M. B. Effects of round-to-nearest and stochastic rounding in the numerical solution of the heat equation in low precision. *ArXiv201016225 Cs Math* (2020).
43. Fasi, M. & Mikaitis, M. Algorithms for stochastically rounded elementary arithmetic operations in IEEE 754 floating-point arithmetic. *IEEE Trans. Emerg. Top. Comput.* 1–1. <https://doi.org/10.1109/TETC.2021.3069165> (2021).
44. Higham, N. J. The accuracy of floating point summation. *SIAM J. Sci. Comput.* **14**, 783–799 (1993).
45. Klöwer, M. StochasticRounding.jl: Up or down, or maybe both? (v0.6.3). *Zenodo*. <https://doi.org/10.5281/zenodo.8131795> (2023).
46. May, R. M. Simple mathematical models with very complicated dynamics. *Nature* **261**, 459–467 (1976).
47. Parry, W. On the  $\beta$ -expansions of real numbers. *Acta Math. Acad. Sci. Hung.* **11**, 401–416 (1960).
48. Rényi, A. Representations for real numbers and their ergodic properties. *Acta Math. Acad. Sci. Hung.* **8**, 477–493 (1957).
49. Alzaidi, A. A., Ahmad, M., Doja, M. N., Solami, E. A. & Beg, M. M. S. A new 1D chaotic map and  $\beta$ -hill climbing for generating substitution-boxes. *IEEE Access* **6**, 55405–55418 (2018).
50. Hofbauer, F.  $\beta$ -Shifts have unique maximal measure. *Monatshefte Für Math.* **85**, 189–198 (1978).
51. Matsumoto, M. & Nishimura, T. Mersenne twister: A 623-dimensionally equidistributed uniform pseudo-random number generator. *ACM Trans. Model. Comput. Simul.* **8**, 3–30 (1998).
52. Marsaglia, G. Xorshift RNGs. *J. Stat. Softw.* **8**, 1–6 (2003).
53. Blackman, D. & Vigna, S. Scrambled linear pseudorandom number generators. *ACM Trans. Math. Softw.* **47**, 36:1–36:32 (2021).
54. Lorenz, E. N. Predictability: A problem partly solved. in *Proc. Seminar on predictability* vol. 1 (1996).
55. Lorenz, E. N. & Emanuel, K. A. Optimal sites for supplementary weather observations: simulation with a small model. *J. Atmos. Sci.* **55**, 399–414 (1998).
56. Hatfield, S., Subramanian, A., Palmer, T. & Düben, P. Improving weather forecast skill through reduced-precision data assimilation. *Mon. Weather Rev.* **146**, 49–62 (2017).
57. Schneider, T., Lan, S., Stuart, A. & Teixeira, J. Earth system modeling 2.0: A blueprint for models that learn from observations and targeted high-resolution simulations. *Geophys. Res. Lett.* **44**, 12396–12417 (2017).
58. Rasp, S. Coupled online learning as a way to tackle instabilities and biases in neural network parameterizations: general algorithms and Lorenz 96 case study (v1.0). *Geosci. Model Dev.* **13**, 2185–2196 (2020).
59. Butcher, J. C. Runge–Kutta Methods. in *Numerical Methods for Ordinary Differential Equations* 137–316 (John Wiley & Sons, Ltd, 2008). <https://doi.org/10.1002/9780470753767.ch3>.
60. Klöwer, M. Lorenz96.jl: A type-flexible Lorenz 1996 model (v0.3.0). *Zenodo*. <https://doi.org/10.5281/zenodo.5121430> (2021).
61. Klöwer, M., Hatfield, S., Croci, M., Düben, P. D. & Palmer, T. N. Fluid simulations accelerated with 16 bit: Approaching 4x speedup on A64FX by squeezing ShallowWaters.jl into Float16. *J. Adv. Model. Earth Syst.* **14**, e2021MS002684. <https://doi.org/10.1029/2021MS002684> (2021).
62. Rasp, S. Lorenz '96 is too easy! Machine learning research needs a more realistic toy model. *Medium* <https://towardsdatascience.com/lorenz-96-is-too-easy-machine-learning-research-needs-a-more-realistic-toy-model-6add938f6cc0> (2020).
63. Villani, C. *Topics in Optimal Transportation*. (American Mathematical Soc., 2003).
64. Klöwer, M. & Paxton, E. A. BernoulliMap: Chaos in one variable. *Zenodo* <https://doi.org/10.5281/zenodo.8138508> (2023).

## Acknowledgements

MK thanks the Natural Environmental Research Council NERC for funding under grant number NE/L002612/1. MK, EAP and TNP gratefully acknowledge funding from the European Research Council under the European Union's Horizon 2020 research and innovation programme for the ITHACA grant (no. 741112). PVC is grateful for funding from European Union Horizon 2020 Research and Innovation Programme grant number 800925 and UK EPSRC grant number EP/R029598/1. We are also grateful for discussions with Bruce Boghosian, Wouter Edeling and Chiara Maiocchi and for the reviews of two anonymous reviewers.

## Author contributions

Conceptualization: M.K., P.V.C., E.A.P. Data curation: M.K. Formal Analysis: M.K. Methodology: M.K., E.A.P. Visualisation: M.K. Writing—original draft: M.K. Writing—review & editing: M.K., P.V.C., E.A.P., T.N.P.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-023-37004-4>.

**Correspondence** and requests for materials should be addressed to M.K.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023