

1 **ROUNDING ERROR USING LOW PRECISION APPROXIMATE**
2 **RANDOM VARIABLES**

3 MICHAEL GILES* AND OLIVER SHERIDAN-METHVEN†

4 **Abstract.** For numerical approximations to stochastic differential equations using the Euler-
5 Maruyama scheme, we propose incorporating approximate random variables computed using low
6 precisions, such as single and half precision. We propose and justify a model for the rounding error
7 incurred, and produce an average case error bound for two and four way differences, appropriate
8 for regular and nested multilevel Monte Carlo estimations. Our rounding error model recovers and
9 extends the statistical model by Arciniega and Allen [1], while bounding the size systematic and
10 biased rounding errors are permitted to be. By considering the variance structure of multilevel Monte
11 Carlo correction terms in various precisions with and without a Kahan compensated summation, we
12 compute the potential speed ups offered from the various precisions. We find single precision offers
13 the potential for approximate speed improvements by a factor of 7 across a wide span of discretisation
14 levels. Half precision offers comparable improvements for several levels of coarse simulations, and
15 even offers improvements by a factor of 10–12 for the very coarsest few levels, which are likely to
16 dominate higher order methods such as the Milstein scheme.

17 **Key words.** approximations, random variables, inverse cumulative distribution functions, ran-
18 dom number generation, finite precision, half precision, floating point, rounding error, multilevel
19 Monte Carlo, the Euler-Maruyama scheme, the Milstein scheme, Kahan compensated summation,
20 and high performance computing.

21 **AMS subject classifications.** 65G50, 65C10, 41A10, 65C05, 65Y20, 60H35, 65B10, 65L70,
22 34M30, 97N20, and 65C30.

23 **1. Introduction.** Rounding error has long been a source of scientific interest and
24 frustration. For a large fraction of the scientific community, the effects of rounding
25 error are usually negligible and of little or no consequence. However, for an appreciable
26 portion of the community, especially those pushing computer hardwares and numerical
27 algorithms to the fastest speeds achievable, rounding error can present a considerable
28 hurdle to the achievable fidelity and speed.

29 The example *par excellence* of rounding error in scientific computing is in sum-
30 mation operations, frequent in linear algebra applications and to a lesser extent some
31 statistical applications. Calculating the scalar product between two vectors, and thus
32 also calculating vector and matrix multiplications, involves (among other things) sum-
33 ming a list of numbers. Being able to accurately sum a list of numbers has long been
34 under the attention of mathematicians and computer scientists [21, 29, 33, 50], and its
35 importance in scientific computing cannot be understated. When the numbers being
36 summed are ill conditioned, the impact of rounding error grows with the problem’s
37 size, and can quickly nullify even the simplest of calculations, and thus high accu-
38 racy summation algorithms are frequently required. Outside of linear algebra, the
39 accuracy of gradients and sensitivity estimates from finite difference methods and nu-
40 merical differentiation is capped by the maximum available precision due to rounding
41 error. Lastly, for the numerical solution of differential equations by simulation meth-
42 ods (deterministic or stochastic), iterative methods such as the Euler scheme incur
43 rounding errors which get worse as the simulation’s discretisation becomes finer.

44 To combat the effects of rounding error, there are two particularly common ap-
45 proaches. The first is to try and bypass the issue by simply working in a higher pre-
46 cision, typically at the cost of computational speed. Historically this has motivated

*mike.giles@maths.ox.ac.uk

†oliver.sheridan-methven@hotmail.co.uk

47 the introduction of double precision, extended double precision, and even quadruple
48 precision data types. Similarly, several software libraries offer arbitrary levels of pre-
49 cision, such as: the mpmath [28] Python library, the GNU multiple precision (GMP)
50 arithmetic C/C++ library [20], and the GNU multiple precision binary floating point
51 with correct rounding (MPFR) C library [12]. In a similar vein, hardware providers
52 have also focussed efforts on improved floating point accuracy and reproducibility [6].

53 The second approach to reduce the influence of rounding error is to try and
54 compensate and correct against it. For summations, the best known approach is the
55 Kahan compensated summation [29], although other compensation procedures have
56 also been introduced and well explored [2, 10, 31, 33, 34, 38, 40, 41]. These proceed by
57 inferring an estimate for the rounding error introduced at each stage of the summation,
58 and then discount for this in the subsequent summations, thus compensating for the
59 rounding error.

60 There is a third means of circumventing rounding error, which is more mathemat-
61 ical in its nature, which looks to extrapolate accurate answers from less accurate ap-
62 proximations. The best example of this is Richardson extrapolation [35, 44]. Without
63 this technique, approximation schemes would need to go to such fine discretisations
64 that rounding error would be significant, whereas using Richardson extrapolation is
65 one possible technique to avoid encountering rounding error. However, the use of
66 Richardson extrapolation is very problem specific, and whilst it is a very powerful
67 mathematical technique, it does not readily present itself as a general purpose com-
68 putational tool for avoiding rounding error in most circumstances.

69 Both high precision libraries and rounding error compensation schemes have his-
70 torically been reserved for specialised applications seeking extraordinarily high accu-
71 racy, and closer to the edges of most scientific computing applications. However, in
72 more recent years there has been an increased demand for ever lower precisions and
73 data types, such as the IEEE half precision float [26], and the more recent “brain
74 float” [7, 30]. The large driving force behind these is the increasingly popular demand
75 in machine learning applications, where the underlying data are very noisy and impre-
76 cise, and smaller data types are preferable for faster access, storage, and computation
77 on the latest CPU, GPU, and TPU hardware. Furthermore, with the greater desire
78 for increased parallelisation on vector hardware and reduced precision calculations,
79 lower precision data types are gaining considerable momentum and traction.

80 To understand the nature of the nett rounding error arising in calculations, there
81 have been two fronts of development. The first has been defining the precise rounding
82 modes and data types used in scientific calculations. To ensure floating point calcula-
83 tions were standardised, the famous IEEE 754 standard for floating point arithmetic
84 was introduced [25], and is now the *de facto* industry standard. This entails addi-
85 tion, subtraction, multiplication, division, and square roots all producing exact results
86 with respect to the appropriate rounding mode [51, page 15]. Similarly, the rounding
87 modes a computer uses to round floating point values are standardised, with “round
88 to nearest even” typically being the default mode. Of course, while floating point
89 arithmetic may be well defined, there are several difficulties and nuances, as discussed
90 by Goldberg [19].

91 The second front has been with the mathematical modelling of rounding errors.
92 While the IEEE 754 standard specifies the hardware’s behaviour, this does not readily
93 give insight into the behaviour of the emergent nett rounding error. Describing the
94 nett effect that results during calculations has received much mathematical attention
95 [22, 24, 53, 54, 55], and one of the best overviews is by Higham [22], who analyses
96 the standard model for deterministic rounding error [22, 2.2, (2.4)]. Furthermore,

97 in recent years there has been a piqued interest in stochastic rounding modes and
98 associated error models, with prominent recent work by Higham and Mary [23] and
99 Ipsen and Zhou [27] performing probabilistic error analyses, which frequently give
100 tighter and more realistic error bounds than the worst case deterministic scenarios.
101 In the setting of analysing partial differential equations using stochastic rounding
102 there is the recent work by Croci and Giles [9].

103 With this wealth of attention from academia and industry, the work we present
104 produces a rigorous model for the rounding error incurred during the numerical simula-
105 tion of stochastic differential equations. Typical treatments of numerical methods for
106 such stochastic differential equations assume no rounding error occurs or is otherwise
107 negligible [32, 9.3, page 316] [18]. The earliest work to compensate for rounding error
108 in simulations for ordinary differential equations appears to be by Vitasek [52]. In the
109 setting of stochastic differential equations, the most relevant works are by Arciniega
110 and Allen [1] and Omland [42]. Arciniega and Allen [1] present an *ad hoc* statistically
111 motivated model for the rounding error which occurs in the Euler-Maruyama scheme,
112 giving an average case bound for the overall rounding error. Omland [42] takes a more
113 rigorous approach, closer to a first principles model, starting with the floating point
114 rounding modes and standard error model by Higham [22], and produces a worst case
115 bound for the error in the Euler-Maruyama scheme [42, theorem 4.8]. Related to
116 modelling rounding error within stochastic simulations is work by Delattre and Jacod
117 [11] and Rosenbaum [45], statistically modelling the effects of simple rounding mod-
118 els on statistical inference problems driven by underlying stochastic processes. They
119 analyse the convergence of empirical statistical estimators for functionals of certain
120 stochastic processes typical in finance, estimating parameters occurring in diffusion
121 coefficients using empirical samples.

122 The contribution of this work will be to produce a heuristic model for the rounding
123 error in a similar manner to Arciniega and Allen [1]. However, our model will be much
124 more rigorously justified by a detailed inspection of the dominant rounding errors
125 arising in the Euler-Maruyama scheme. Furthermore, we will show that there are two
126 primary sources of error which contribute to the nett error. The first is a zero mean
127 process, similar to that described by Arciniega and Allen [1], and any deviation from
128 being zero mean is comprehensively analysed, quantified, and shown to be negligible
129 compared to the other error processes occurring. The second is a possibly non zero
130 mean systematic error term omitted by Arciniega and Allen [1]. This second process
131 is of a much smaller size than the first, but due to its non zero mean nature, we show
132 that its nett contribution will grow at the same rate as that arising from the zero
133 mean process.

134 The significance of this new model is that it both rigorously grounds the leading
135 order error process as being zero mean and now additionally quantifies the permis-
136 sible size of lower order systematic rounding errors in the Euler-Maruyama scheme.
137 Furthermore, our model is amenable to incorporation within the nested multilevel
138 Monte Carlo scheme utilising approximate random variables developed by Giles and
139 Sheridan-Methven [15, 46, 47, 49]. Although work has been done by Brugger et al. [5]
140 and Omland et al. [43] on constructing multilevel Monte Carlo schemes in the pres-
141 ence of rounding error, our model directly facilitates a treatment jointly allowing for
142 approximate random variables and low precision calculations, only correctly handled
143 by a nested multilevel Monte Carlo scheme.

144 A secondary contribution of this work will also be to demonstrate the applicability
145 of a Kahan compensated summation within the Euler-Maruyama scheme, an extension
146 of the similar idea by Vitasek [52] in the setting of ordinary differential equations. The

147 inclusion of Kahan compensated summation we show considerably increases the range
148 of applicability of half precision to many more discretisation levels beyond just the
149 coarsest few.

150 The combination of these two contributions should provide considerable interest
151 to both theoreticians, but crucially also practitioners working in high performance
152 scientific computing. Presenting concrete mathematical underpinnings for modelling
153 the effect of rounding errors, our work newly describes the impact of both leading
154 order zero mean and second order systematic error processes on the nett rounding
155 error that arises. Being able to bound the permissible size of systematic processes is
156 a novel result now available to those wishing to model or construct further numerical
157 schemes building off this work. Additionally, integrating these results in the nested
158 multilevel Monte Carlo framework and demonstrating the utility of half precision and
159 Kahan compensated summation opens up our work to practitioners. The empirical
160 results demonstrate the substantial benefits to be had by adopting our framework
161 and heavily engaging with low precision floating point formats such as half precision,
162 whilst our analytic result give the necessary reassurance and confidence that accuracy
163 can be all the while maintained, and need not be sacrificed.

164 Section 2 overviews the numerical solution of stochastic differential equations,
165 providing the primary context and setting of our work, presenting our model for the
166 leading order error process arising in the Euler-Maruyama scheme. Section 3 will
167 showcase how our model can be incorporated into a multilevel Monte Carlo frame-
168 work, demonstrating practical applications of the model and highlighting the temporal
169 savings that can be expected using low precisions. Section 4 presents the conclusions
170 from this work.

171 **2. Numerical solutions to stochastic differential equations.** There are
172 various settings appropriate for analysing the effects of rounding error, and the nu-
173 merical solutions of stochastic differential equations is one particularly important
174 setting. Frequently the terminal solution X_T of the stochastic differential equation
175 $dX_t = a(t, X_t) dt + b(t, X_t) dW_t$ needs to be approximated for given drift and diffusion
176 processes a and b . To achieve this, whole path approximations $\hat{X}_t \approx X_t$ for $t \in [0, T]$
177 are produced, where the most popular methods are the Euler-Maruyama and Milstein
178 schemes. For a thorough detailing see Kloeden and Platen [32] and Glasserman [18].
179 The approximations simulate the process over N time steps of size $\Delta t \equiv \delta := \frac{T}{N}$, where
180 the update at the n -th iteration at time $t_n := n\delta$ requires a Wiener process increment
181 ΔW . The usual numerical schemes use a standard Gaussian random variable Z_n to
182 simulate from this process, where $\Delta W_n := \sqrt{\delta}Z_n$, and these Z_n are independently
183 and identically distributed.

184 To ensure the stochastic process has a unique strong solution and the Euler-
185 Maruyama scheme converges, we assume the standard assumptions from Kloeden and
186 Platen [32, 4.5, 10.2], which are that: a and b are jointly Lebesgue measurable, spa-
187 tially Lipschitz continuous, have linear spatial growth, $\frac{1}{2}$ -Hölder temporal continuity
188 with linear spatial growth, and that X has a measurable initial condition.

189 The rounding error model we will later introduce will be applicable to the original
190 Euler-Maruyama scheme, but also to a modified version utilising high speed and low
191 fidelity “approximate random variables”. Consequently, we will now review approx-
192 imate random variables and detail the associated requirements needed for our error
193 model.

194 **2.1. Approximate random variables.** Unfortunately, sampling from the Gauss-
195 ian distribution is expensive, and so there have been several approaches to bypass this

196 cost. Most of these look to substitute the exact Gaussian increment Z_n with another
 197 random variable \tilde{Z}_n with similar statistics. For clarity and consistency with Giles
 198 and Sheridan-Methven [15, 16, 46], we call these substitutes *approximate random*
 199 *variables*, and the originals as *exact random variables*. The most well known is to use
 200 Rademacher random variables, producing what’s known as the weak Euler-Maruyama
 201 scheme [32, page XXXII], where the Rademacher random variables have the desired
 202 mean. More advanced methods include more generalised moment matching proce-
 203 dures, as discussed by Muller [39], and piecewise polynomial approximations and
 204 generalised approximate random variables by Giles and Sheridan-Methven [15, 46].

205 We briefly pause to remark that the aforementioned methods for generating
 206 Gaussian random variables can all be viewed as utilising approximations to the Gauss-
 207 ian distribution’s inverse cumulative distribution function in the inverse transform
 208 method [15, 18]. However, there are alternative methods for generating Gaussian
 209 random variables, including the Ziggurat [37], Box-Muller [4], and Marsaglia [36]
 210 methods, each being extensively used by practitioners. For our purposes, we will
 211 be assuming the exact random variables are produced using the inverse transform
 212 method, and that the approximate random variables are generated using the piece-
 213 wise polynomial approximations by Giles and Sheridan-Methven [16]. There are two
 214 main reasons for this. The first is that we will be able to utilise the body of analytic
 215 results derived by Giles and Sheridan-Methven [15, 16] concerning both the approx-
 216 imate random variables, and their subsequent use in nested multilevel Monte Carlo
 217 schemes. The second is that using the approximate random variables produced us-
 218 ing [16] have been demonstrated to be substantially faster than their exact random
 219 variable counterparts on modern vectorised hardware [16]. This speed improvement
 220 is by design, and scales to smaller data types and wider vector lengths, whereas the
 221 alternatives methods do not, and their limitations are discussed extensively in [16,
 222 section 3.1.1], which we don’t repeat here.

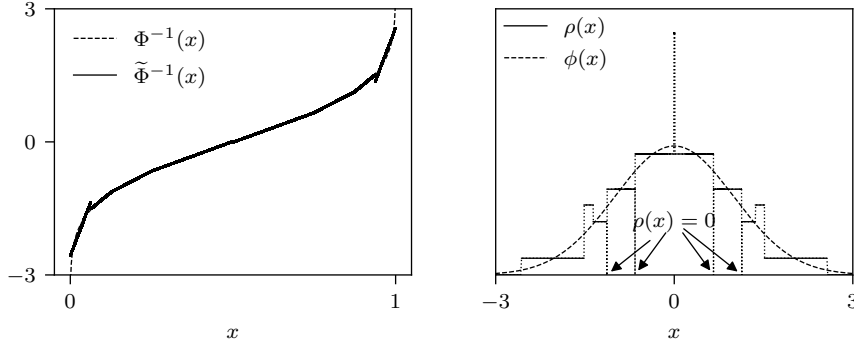
223 The Euler-Maruyama schemes using the exact Gaussian random variables Z_n give
 224 rise to the approximation \hat{X} , and the approximate random variables \tilde{Z}_n produce \tilde{X} ,
 225 where the schemes are respectively

$$226 \quad \hat{X}_{n+1} = \hat{X}_n + a(t_n, \hat{X}_n)\delta + b(t_n, \hat{X}_n)\sqrt{\delta}Z_n$$

228 and

$$229 \quad \tilde{X}_{n+1} = \tilde{X}_n + a(t_n, \tilde{X}_n)\delta + b(t_n, \tilde{X}_n)\sqrt{\delta}\tilde{Z}_n.$$

231 One of the approximations we will utilise later is the piecewise linear approxi-
 232 mation by Giles and Sheridan-Methven [16]. This generates approximate Gaussian
 233 random variables by the inverse transform method [18] using a piecewise linear approx-
 234 imation $\tilde{\Phi}^{-1}$ to the Gaussian distribution’s inverse cumulative distribution function
 235 Φ^{-1} . The exact construction is detailed by Giles and Sheridan-Methven [16], although
 236 an example approximation is demonstrated in Figure 2.1. A piecewise linear approx-
 237 imation using 8 intervals is shown in Figure 2.1(a), and the resultant probability
 238 density function of the approximation ρ is shown in Figure 2.1(b). The piecewise lin-
 239 ear approximation is non-monotonic with small discontinuities at the interval bound-
 240 aries. Consequently, the resulting probability density function has compact support,
 241 although there are a few resulting tiny non-obvious inaccessible domain intervals with
 242 zero measure, as indicated. The rounding error model we will later propose will hold
 243 for both exact Gaussian random variables, and also certain classes of approximate
 244 Gaussian random variables, for which we require that they satisfy Assumption 2.1.



(a) A piecewise linear approximation using 8 intervals.

(b) The resultant probability density function with non-obvious intervals of zero measure indicated.

FIG. 2.1. A piecewise linear approximation of the Gaussian distribution's inverse cumulative distribution function, and the resultant probability density function.

245 ASSUMPTION 2.1. Let any approximate Gaussian random variables \tilde{Z} be zero
 246 mean, uniformly bounded, and have finite variance $\mathbb{V}(\tilde{Z}) = O(1)$, and also have all
 247 higher order moments be finite. Furthermore, let there exist a corresponding proba-
 248 bility density function ρ such that $\mathbb{P}(\tilde{Z} \in [z, z + dz]) = \rho(z) dz$. Let ρ be bounded by
 249 K such that $\rho \leq K < \infty$ and be smooth almost everywhere such that $\rho \in C^\infty(\mathbb{R} \setminus \mathcal{M})$,
 250 where \mathcal{M} is a finite set of M points where ρ is discontinuous. Lastly, let ρ decay suffi-
 251 ciently fast such that for any finite constant α we have $\sum_{k=-\infty}^{\infty} 2^{2k} \max_{y' \in [2^k, 2^{k+1}]} \rho(y' -$
 252 $\alpha) < \infty$.

253 LEMMA 2.1. The exact Gaussian distribution satisfies Assumption 2.1.

254 *Proof.* We immediately have that the Gaussian distribution is zero mean and has
 255 unit variance, and is uniformly bounded [3, appendix C.2]. The probability density
 256 function for the Gaussian distribution is $\rho \equiv \phi$ where $\phi(z) := \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}z^2)$, which is
 257 $C^\infty(\mathbb{R})$ and maximal at zero where $\rho \leq \phi(0) = \frac{1}{\sqrt{2\pi}}$. To show $\sum_{k=-\infty}^{\infty} 2^{2k} \max_{y' \in [2^k, 2^{k+1}]} \rho(y' -$
 258 $\alpha) < \infty$ we consider separately the cases when $\alpha > 0$ and $\alpha \leq 0$. When $\alpha > 0$, we note
 259 that the summand is increasing as k increases from $-\infty$, and ρ will be maximal for
 260 the one index k^* where $\alpha \in [2^{k^*}, 2^{k^*+1}]$. The maximal value will therein be obtained
 261 for $y' = \alpha$, and thereafter for indices $k > k^*$ the ρ term will be decreasing. Thus we
 262 can approximate the possibly divergent part of the summation by the integral

$$\begin{aligned}
 263 \quad \sum_{k=k^*+1}^{\infty} 2^{2k} \max_{y' \in [2^k, 2^{k+1}]} \rho(y' - \alpha) &= \sum_{k=k^*+1}^{\infty} 2^{2k} \rho(2^k - \alpha) \\
 264 &\approx \int_{2^{k^*+1}}^{\infty} x^2 \rho(x - \alpha) dx \\
 265 &\leq \int_{-\infty}^{\infty} x^2 \rho(x - \alpha) dx \\
 266 &< \infty.
 \end{aligned}$$

268 The summation can be strictly bounded from above by using similar integrals, and

269 thus is not divergent. The partial summation from $k = -\infty$ up to $k = k^*$ is immedi-
 270 ately finite as ρ is bounded.

271 For the case $\alpha \leq 0$, in each $[2^k, 2^{k+1}]$ interval the ρ is maximal at $y' = 2^k$, and
 272 decreases as k increases. The relevant summation in this case can also be bounded
 273 identically and shown to be non divergent. \square

274 LEMMA 2.2. *The approximate Gaussian distribution resulting from the piecewise*
 275 *linear approximation by Giles and Sheridan-Methven [16] satisfies Assumption 2.1.*

276 *Proof.* For a finite number of approximation intervals, the probability density
 277 function is symmetric, uniformly bounded, and has compact support and thus finite
 278 variance. The number of discontinuities is finite, and as ρ has compact support it
 279 immediately satisfies the summation bound from Assumption 2.1. \square

280 LEMMA 2.3. *The approximate Gaussian distribution resulting from the piecewise*
 281 *cubic approximation by Giles and Sheridan-Methven [16] satisfies Assumption 2.1.*

282 *Proof.* The proof follows identically to the proof of Lemma 2.2. \square

283 The motivation for introducing these approximate random variables was to in-
 284 crease simulation speed. However, for the ultimate speed improvements, it is desir-
 285 able to both switch to approximate random variables, and simultaneously decrease
 286 the arithmetic precision used, giving a twofold speed improvement. Reducing the
 287 precision alone has been explored with applications to field programmable gate ar-
 288 rays [5, 8, 42, 43], as have multilevel Monte Carlo schemes using varying fidelities
 289 of approximate random variables [39]. However, performing both simultaneously is
 290 touched upon by Giles et al. [17], although using a quite restrictive truncated uniform
 291 random bit Monte Carlo algorithm, and extensions to more general approximation
 292 schemes without varying the precision is done by Giles and Sheridan-Methven [15].
 293 However, the work by Giles et al. [17] is primarily a cost most, and does not model
 294 the effect of rounding error. Thus, while our contribution is an extension of these
 295 works, it is an important and novel demonstration and vindication of the utility of
 296 low precisions with approximate random variables.

297 In order to describe the effect of rounding error resulting from the Euler-Maruyama
 298 scheme, the two most prominent works are by Arciniega and Allen [1] and Omland
 299 [42], which are models for the average and worst case errors respectively. Arcin-
 300 iega and Allen [1] provide an *ad hoc* statistical model and analysis for the round-
 301 ing error arising from finite precision floating point calculations within the Euler-
 302 Maruyama scheme. Denoting the estimate produced when working in finite precision
 303 as \bar{X} , they propose that at the n -th iteration, all of the composite floating point
 304 arithmetic in the Euler-Maruyama update culminates in an additive error ε_n where
 305 $\bar{X}_{n+1} = \bar{X}_n + a(t_n, \bar{X}_n)\delta + b(t_n, \bar{X}_n)\sqrt{\delta}Z_n + \varepsilon_n$. This error is assumed to follow a
 306 Gaussian distribution, be zero mean, and have a variance $\mathbb{V}(\varepsilon) \leq C\varrho^2$, where ϱ is
 307 the unit roundoff and C is some constant. The main result from their analysis [1,
 308 theorem 2.2] is $\mathbb{E}(|\hat{X}_N - \bar{X}_N|^2) \leq CN^2\varrho^2$. The model from Omland [42] uses a more
 309 rigorous finite precision framework. For brevity, letting \oplus and \otimes represent floating
 310 point addition and multiplication respectively, then the model by Omland [42] in ef-
 311 fect considers $\bar{X}_{n+1} = \bar{X}_n \oplus ((a(t_n, \bar{X}_n) \otimes \delta) \oplus ((b(t_n, \bar{X}_n) \otimes \sqrt{\delta}) \otimes Z_n))$ and produces
 312 the worst case bound $\mathbb{E}(|\hat{X}_N - \bar{X}_N|^2) \leq CN^2\varrho^2$ [42, theorem 4.8]. The model we
 313 will present will take the model from Omland [42] as its starting point, but under
 314 assumptions appropriate for the Euler-Maruyama scheme, will ultimately reduce to
 315 a model closer resembling that by Arciniega and Allen [1]. Our model will also allow
 316 for systematic rounding effects, and bounds how large these are permitted to be and

317 have the model still be applicable.

318 **2.2. A leading order error model.** Starting with the more fundamentally
 319 rooted model from Omland [42], we expect to recover, under appropriate assumptions,
 320 the more statistically motivated model from Arciniega and Allen [1]. To this end we
 321 look to expand the model from Omland [42] and see the effects of arithmetic roundoff
 322 within the Euler-Maruyama update. Before presenting the analysis, we will briefly
 323 introduce a small amount of floating point notation.

324 Numbers used in calculations must be stored in finite precision, where we denote
 325 the set of representable numbers as $\overline{\mathbb{R}} \subset \mathbb{R}$, where we introduce the rounding operator
 326 $R: \mathbb{R} \rightarrow \overline{\mathbb{R}}$ which implements the desired rounding mode, which we assume is round to
 327 nearest even. For finite precision binary arithmetic operations \oplus (such as \oplus , \otimes , etc.),
 328 then for two floating point numbers $x, y \in \overline{\mathbb{R}}$, we assume the standard rounding model
 329 by Higham [22, 2.2, (2.4)] [50, page 99, (13.7)] that $x \oplus y = R(x * y) = (x * y)(1 + \varepsilon)$
 330 where $|\varepsilon| \leq \varrho$.

331 For further notational convenience, when we talk about a random variable or
 332 process g being of a typical size h , we usually mean in the sense of absolute size and
 333 strong expectation, where we are really discussing the size of $\mathbb{E}(|g|)$. However, such
 334 explicit notation is overly verbose, so instead we will make regular use of O -notation.
 335 For brevity, we will write statements such as $\mathbb{E}(|g|) = O(h)$ instead as $|g| = O(h)$
 336 or even $g = O(h)$, where the underlying meaning of such statements should be clear
 337 from context.

338 We begin by expressing the finite precision Euler-Maruyama update as $\overline{X}_{n+1} =$
 339 $\overline{X}_n \oplus (A_n \oplus B_n)$, where $A_n := a(t_n, \overline{X}_n) \otimes \delta$ and $B_n := (b(t_n, \overline{X}_n) \otimes \sqrt{\delta}) \otimes \tilde{Z}_n$. For an
 340 appropriately non dimensionalised stochastic process, such that $T = 1$, $X_0 = O(1)$,
 341 $a = O(1)$, $b = O(1)$, then we anticipate $\overline{X}_n = O(1)$, $\tilde{Z}_n = O(1)$, $A_n = O(\delta)$, and
 342 $B_n = O(\sqrt{\delta})$. Similarly, for the precision levels and discretisations we have $\varrho \ll 1$ and
 343 $\delta \ll \sqrt{\delta} \ll 1$. This then gives us the size ordering $\mathbb{E}(|A_n|) \ll \mathbb{E}(|B_n|) \ll \mathbb{E}(|\overline{X}_n|)$.

344 The first addition $A_n \oplus B_n$ will produce an absolute error η'_n where $A_n \oplus B_n =$
 345 $A_n + B_n + \eta'_n$, and η'_n will be of a size comparable with the unit roundoff and the larger
 346 of A_n and B_n , which is B_n . As b is assumed to have linear growth we obtain $|\eta'_n| \sim$
 347 $|B_n \varrho| = O(\sqrt{\delta} \varrho (1 + |\overline{X}_n|))$, and after performing this first floating point addition we
 348 will be left with $\overline{X}_{n+1} = \overline{X}_n \oplus (B_n + A_n + \eta'_n)$, where we have written $B_n + A_n + \eta'_n$
 349 in order of decreasing magnitudes.

350 For the remaining addition operation $\overline{X}_n \oplus \dots$, as we have $X_n = O(1)$ and
 351 $\mathbb{E}(|\overline{X}_n|) \gg \mathbb{E}(|B_n|)$, the nett result from the remaining floating point addition then
 352 is that this will produce a second absolute arithmetic error η_n where $|\eta_n| \sim |\overline{X}_n \varrho| =$
 353 $O(\varrho (1 + |\overline{X}_n|))$, and the Euler-Maruyama update will become $\overline{X}_{n+1} \approx \overline{X}_n + B_n +$
 354 $A_n + \eta_n + \eta'_n$, where again we have written the contributions in order of decreasing
 355 magnitudes. We identify two dominant sources of error. The first is η'_n , arising from
 356 the addition of the drift term to the diffusion term. The second is η_n , arising from the
 357 addition of this sum to the underlying process. We expect $|\eta'_n| = O(\varrho \sqrt{\delta} (1 + |\overline{X}_n|))$
 358 and $|\eta_n| = O(\varrho (1 + |\overline{X}_n|))$, and thus $\mathbb{E}(|\eta_n|) \gg \mathbb{E}(|\eta'_n|)$. If we then allow for the
 359 inclusion of other higher order contributions η'' , we can then see that we expect to
 360 obtain $\overline{X}_{n+1} = \overline{X}_n + B_n + A_n + \eta_n + \eta'_n + \eta''$.

361 The Arciniega and Allen [1] model assumes that only η_n is significant, and makes
 362 the assertion that this can be modelled as a zero mean Gaussian random variable
 363 with a variance only proportional ϱ^2 . In our model, we will more rigorously justify
 364 the zero mean nature, drop the requirement that this exactly follows a Gaussian
 365 distribution, and show that the smaller second order contributions from η' are not

366 necessarily negligible, but can contribute to the nett rounding error at the same rate
 367 as the leading order η process. We propose Model 2.1 as an appropriate model for
 368 the rounding errors arising in the Euler-Maruyama scheme.

369 **MODEL 2.1.** *Let the Euler-Maruyama scheme use random variables \tilde{Z} which sat-*
 370 *isfy Assumption 2.1. The composite effects of rounding error introduce two dominant*
 371 *sources of error, η and η' , where at each step we have*

$$372 \quad \bar{X}_{n+1} = \bar{X}_n + a(t_n, \bar{X}_n)\delta + b(t_n, \bar{X}_n)\sqrt{\delta}\tilde{Z}_n + \eta_n + \eta'_n.$$

373 *The larger of these is $\eta_n = O(\varrho(1 + |\bar{X}_n|))$, which is a martingale increment, and*
 374 *the smaller of these is $\eta'_n = O(\varrho\sqrt{\delta}(1 + |\bar{X}_n|))$, which is a possibly non martingale*
 375 *increment.*

376 Inspecting Model 2.1, the key modelling assumption requiring justification is the
 377 martingale¹ nature of η_n . We already justified in our discussions that $\eta'_n = O(\varrho\sqrt{\delta}(1 +$
 378 $|\bar{X}_n|))$, and so claiming this is a non martingale increment is no further restriction.
 379 It is straightforward to reason that $\mathbb{E}(|\eta_n|) = O(\varrho(1 + \mathbb{E}(|\bar{X}_n|)))$, which if we take
 380 $\mathbb{E}(|\bar{X}_n|) = O(1)$ simplifies to $\mathbb{E}(|\eta_n|) = O(\varrho)$. Thus, to justify η_n being a martingale
 381 increment it is sufficient to reason that $|\mathbb{E}(\eta_n)| \ll \mathbb{E}(|\eta_n|)$, which we achieve through
 382 Lemma 2.4.

383 **LEMMA 2.4.** *Assuming $\mathbb{E}(|\bar{X}_n|) = O(1)$ and the random variables \tilde{Z} satisfy As-*
 384 *sumption 2.1, then under the round to nearest even rounding mode, the leading order*
 385 *rounding error η_n has $\mathbb{E}(\eta_n) = O(\varrho^2)$.*

386 *Proof.* The operation producing η_n is the floating point addition between \bar{X}_n and
 387 $(A_n \oplus B_n)$. Dropping the subscripts for brevity, we denote this by $\alpha \oplus R(\beta)$, where
 388 $\alpha := \bar{X}_n = O(1)$ and $\beta := A_n + B_n = O(\sqrt{\delta})$. The absolute rounding error η is then
 389 given by

$$390 \quad \eta := (\alpha \oplus R(\beta)) - (\alpha + R(\beta))$$

$$391 \quad \equiv (R(\alpha + \beta) - (\alpha + \beta)) - (R(\beta) - \beta) + (R(\alpha + R(\beta)) - R(\alpha + \beta)),$$

393 where we will bound the final three parenthesised differences in turn.

394 Inspecting the first term, we can see this is the absolute error resulting from
 395 rounding the quantity $\alpha + \beta$. Without much loss of generality, as we have assumed
 396 $\alpha = O(1)$, let us suppose $\alpha \in (1, 2)$. Using IEEE floating point representation, the set
 397 of representable numbers $(1, 2) \cap \bar{\mathbb{R}}$ will all be equally spaced. Defining the quantity
 398 $z := \alpha + \beta$, this will fall inside some interval $I_y := [y - \varsigma, y + \varsigma]$, where $y \pm \varsigma$ are
 399 adjacent floating point numbers. Without loss of generality, we assume $y - \varsigma$ is odd
 400 and $y + \varsigma$ is even, where we either have $z < y$ and we round down, or $z \geq y$ and we
 401 round up, as depicted in Figure 2.2(a).

402 We can then evaluate the expectation $\mathbb{E}((R(z) - z)\mathbf{1}_{\{z \in I_y\}})$ where we have

$$403 \quad \mathbb{E}((R(z) - z)\mathbf{1}_{\{z \in I_y\}}) = \mathbb{E}((R(z) - z)\mathbf{1}_{\{z \in [y - \varsigma, y)\}}) + \mathbb{E}((R(z) - z)\mathbf{1}_{\{z \in [y, y + \varsigma)\}}).$$

404 These expectations can be written as integrals, giving

$$405 \quad \mathbb{E}((R(z) - z)\mathbf{1}_{\{z \in I_y\}}) = \int_{y - \varsigma}^y ((y - \varsigma) - z)\mathbb{P}(dz) + \int_y^{y + \varsigma} ((y + \varsigma) - z)\mathbb{P}(dz).$$

¹For readers unfamiliar with martingales (and stochastic processes), cf. [32, 2.3].

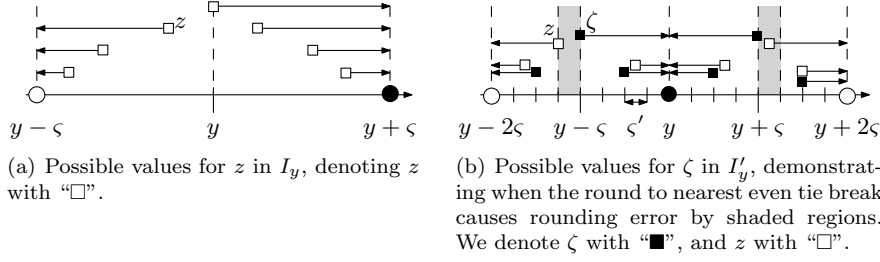


FIG. 2.2. Rounding to the nearest even. We denote even representable values using “●” and odd values using “○”. Arrows show the values rounded to.

406 At time t_n , the variable \bar{X}_n is \mathcal{F}_n -measurable, but \tilde{Z}_n is \mathcal{F}_{n+1} -measurable, and thus z
407 and β will have the same distribution as \tilde{Z} . Denoting the probability density function
408 of β as ρ , which by extension satisfies Assumption 2.1, we obtain

$$409 \quad \mathbb{E}((R(z) - z)\mathbb{1}_{\{z \in I_y\}}) = \int_0^\varsigma (\varsigma - z)(\rho(y - \alpha + z) - \rho(y - \alpha - z)) dz.$$

410 As $\varsigma \ll 1$, if ρ is smooth everywhere in the interval I_y , then we can approximate this
411 using a Taylor series expansion, otherwise we use the bound from Assumption 2.1,
412 obtaining

$$413 \quad |\mathbb{E}((R(z) - z)\mathbb{1}_{\{z \in I_y\}})| \leq \begin{cases} C\varsigma^3|\rho'(y - \alpha)| + O(\varsigma^5|\rho'''(y - \alpha)|) & \text{if } I_y \cap \mathcal{M} = \emptyset \\ C\varsigma^2K & \text{if } I_y \cap \mathcal{M} \neq \emptyset, \end{cases}$$

414 for some constant C .

415 Using this expectation in the law of total expectation, then in the limits $I_y \rightarrow dI_y$
416 we obtain

$$417 \quad \mathbb{E}(R(z) - z) = \mathbb{E}(\mathbb{E}(R(z) - z \mid z \in I_y))$$

$$418 \quad \approx \int_{\mathbb{P}(I_y) > 0} \frac{\mathbb{E}((R(z) - z)\mathbb{1}_{\{z \in I_y\}})}{\mathbb{P}(I_y)} \mathbb{P}(dI_y)$$

$$419 \quad \approx \int \frac{\mathbb{E}((R(z) - z)\mathbb{1}_{\{z \in I_y\}})}{2\varsigma} dy,$$

$$420$$

421 where in the last approximation we used $\mathbb{P}(I_y) \approx \rho(y - \alpha)2\varsigma$ and $\mathbb{P}(dI_y) \approx \rho(y - \alpha) dy$.
422 In the limit $\delta \ll 1$ then we can approximate our integration domain as $y \in [0, \infty)$,
423 giving

$$424 \quad |\mathbb{E}(R(z) - z)| \leq \left| \int_0^\infty C\varsigma^2\rho'(y - \alpha)\mathbb{1}_{\{I_y \cap \mathcal{M} = \emptyset\}} dy \right| + \int_0^\infty C\varsigma K\mathbb{1}_{\{I_y \cap \mathcal{M} \neq \emptyset\}} dy.$$

425 In IEEE representation, ς is a constant between powers of two, and thus $2\varsigma = \varrho 2^k$ for
426 $k \in \mathbb{Z}$. Consequently, the first integral is readily decomposed into sub intervals where
427 ς is a constant. For the second integral we let I_m^* denote the interval containing the
428 discontinuity at position $m \in \mathcal{M}$. These discontinuities occur at the corresponding k

429 values k_m such that $I_m^* \subset [2^{k_m}, 2^{k_m+1}]$, and so we obtain

$$430 \quad |\mathbb{E}(R(z) - z)| \leq C \left| \sum_{k=-\infty}^{\infty} \varrho^2 2^{2k} \int_{2^k}^{2^{k+1}} \rho'(y - \alpha) \mathbb{1}_{\{I_y \cap \mathcal{M} = \emptyset\}} dy \right| + CK \sum_{m \in \mathcal{M}} \varsigma \int_{I_m^*} dy.$$

431 The first integral can largely be evaluated exactly. The integration domain will contain
 432 at most M intervals containing singularities, and thus $M + 1$ subintervals of the form
 433 $\int_a^b \rho'(y - \alpha) dy$ where ρ' is continuous in the domain $[a, b]$, and thus $\int_a^b \rho'(y - \alpha) dy =$
 434 $[\rho(y - \alpha)]_a^b$. We then use the bound $|\int_a^b \rho'(y - \alpha) dy| \leq |\rho(a - \alpha)| + |\rho(b - \alpha)| \leq$
 435 $2 \max_{y' \in [a, b]} \rho(y' - \alpha)$. For the second integral we use $\int_{I_m^*} dy = \varsigma$ to give

$$436 \quad |\mathbb{E}(R(z) - z)| \leq C \varrho^2 \sum_{k=-\infty}^{\infty} 2^{2k} (M+1) \max_{y' \in [2^k, 2^{k+1}]} \rho(y' - \alpha) + CK \varrho^2 \sum_{m \in \mathcal{M}} 2^{2k_m} = O(\varrho^2),$$

437 where in the last equality we used Assumption 2.1. This shows the desired $O(\varrho^2)$ holds
 438 for the first parenthesised error constituting η . The same line of reasoning similarly
 439 holds the for $R(\beta) - \beta$ term also constituting η (akin to taking $\alpha \rightarrow 0$), arriving at
 440 an identical limiting bound.

441 It remains to bound the final $R(\alpha + R(\beta)) - R(\alpha + \beta)$ term constituting η . Unlike
 442 the previous two terms, we will see that this term only contributes a rounding error
 443 when the round to nearest even tie break rule is required, and in most scenarios
 444 $\alpha + R(\beta)$ and $\alpha + \beta$ will round to the the same number. To tackle this final term we
 445 introduce the slightly larger interval $I'_y := [y - 2\varsigma, y + 2\varsigma]$, where $y - 2\varsigma$, y , and $y + 2\varsigma$
 446 are all representable and adjacent. Without loss of generality we assume y is even and
 447 $y \pm 2\varsigma$ are odd. Given $\beta = O(\sqrt{\delta})$ and $\alpha = O(1)$, we know that $|\beta - R(\beta)| \leq \varsigma'$ where
 448 $\varsigma' \ll \varsigma$, and thus β is rounded first on a much finer granularity than $\alpha + \beta$. Keeping
 449 our definition $z := \alpha + \beta$ and introducing $\zeta := \alpha + R(\beta)$, then the discrete set of values
 450 ζ can take has a much finer granularity than the three representable numbers in I'_y ,
 451 namely $\zeta \in \{y \pm n\varsigma'\} \cap I'_y$ for integers $n \in \mathbb{N}$. We display the set of values ζ can take
 452 in I'_y in Figure 2.2(b), where we demonstrate several possible rounding scenarios.

453 Inspecting Figure 2.2(b), we can see that in most situations ζ and z round to the
 454 same number. These only round to different numbers when ζ lies on a tie break value
 455 and z takes a different value that is rounded to an odd number, as indicated by the
 456 shaded regions in Figure 2.2(b). Thus, for the final term we have the expectation

$$457 \quad E((R(\zeta) - R(z)) \mathbb{1}_{\{z \in I'_y\}}) = E((R(\zeta) - R(z)) \mathbb{1}_{\{z \in [y - \varsigma - \varsigma', y - \varsigma]\}} \mathbb{1}_{\{\zeta = y - \varsigma\}}) \\ 458 \quad + E((R(\zeta) - R(z)) \mathbb{1}_{\{z \in (y + \varsigma, y + \varsigma + \varsigma']\}} \mathbb{1}_{\{\zeta = y + \varsigma\}}).$$

460 The first expectation will round $\zeta \rightarrow y$ and $z \rightarrow y - 2\varsigma$, giving a nett rounding error
 461 of 2ς , and the second expectation will round $\zeta \rightarrow y$ and $z \rightarrow y + 2\varsigma$ giving a rounding
 462 error of -2ς . Thus, expressing these expectations as integrals we obtain

$$463 \quad E((R(\zeta) - R(z)) \mathbb{1}_{\{z \in I'_y\}}) = \int_{y - \varsigma - \varsigma'}^{y - \varsigma} 2\varsigma \mathbb{P}(dz) + \int_{y + \varsigma}^{y + \varsigma + \varsigma'} -2\varsigma \mathbb{P}(dz) \\ 464 \quad = 2\varsigma \int_0^{\varsigma'} (\rho(z + y - \varsigma - \varsigma' - \bar{a}) - \rho(z + y + \varsigma - \bar{a})) dz. \\ 465$$

466 For this final integral expression we can again either perform a Taylor series expansion

467 or use our bounds from Assumption 2.1 to obtain to leading order

$$468 \quad |E((R(\zeta) - R(z))\mathbb{1}_{\{z \in I'_y\}})| \leq \begin{cases} C\zeta(\zeta\zeta' + (\zeta')^2)|\rho'(y - \alpha)| & \text{if } I'_y \cap \mathcal{M} = \emptyset \\ C\zeta\zeta'K & \text{if } I'_y \cap \mathcal{M} \neq \emptyset. \end{cases}$$

469 As $\zeta' \ll \zeta$, we see that this bound is equivalent to that found earlier for the other two
 470 terms constituting η . Again, by using the law of total expectation and the same steps
 471 as before we obtain the same $O(\varrho^2)$ bound. Combining the three bounds completes
 472 the proof. \square

473 While Model 2.1 is justified by Lemma 2.4, we can appreciate that the proof of
 474 Lemma 2.4 makes use of a few *ad hoc* approximations and limiting cases. However, our
 475 ultimate aim is only to justify our model, rather than derive it, so such conveniences
 476 are permissible. This serves to illustrate from first principles why the leading order
 477 error term η is effectively zero mean. Overall then, Lemma 2.4 provides a much more
 478 rigorous justification of the zero mean nature of the leading order error than was
 479 simply asserted in the Arciniega and Allen [1] model.

480 We now have $\mathbb{E}(|\eta_n|) = O(\varrho)$, $|\mathbb{E}(\eta_n)| = O(\varrho^2)$, and $\mathbb{E}(|\eta'_n|) = O(\varrho\sqrt{\delta})$. Conse-
 481 quently, for any reasonable discretisation where $\sqrt{\delta} \gg \varrho$, any non-martingale behav-
 482 iour exhibited by the η_n process is negligible compared to the η'_n process, and instead
 483 constitutes part of the η''_n process mentioned earlier.

484 To illustrate how the bound for the final nett rounding error is produced from our
 485 model for the incremental rounding error, we recall a convenient lemma from Giles
 486 and Sheridan-Methven [15, lemma 4.3] [49, lemma 5.2.3], which we present without
 487 proof as Lemma 2.5.

488 LEMMA 2.5. *Suppose for a process \mathcal{E}_n we have $\mathcal{E}_{n+1} = \mathcal{E}_n + \delta\mathcal{A}_n + \sqrt{\delta}\tilde{Z}_n\mathcal{B}_n +$
 489 $\Xi_n + \Theta_n$, using a discretisation interval δ . We assert $\mathcal{E}_0 = 0$ almost surely, \tilde{Z}_n are
 490 i.i.d. zero mean random variables with all finite moments bounded, and \mathcal{A}_n and \mathcal{B}_n
 491 are \mathcal{F}_{t_n} -adapted with $|\mathcal{A}_n| \leq L_A|\mathcal{E}_n|$ and $|\mathcal{B}_n| \leq L_B|\mathcal{E}_n|$ for some strictly positive and
 492 finite constants L_A and L_B . The process Ξ is a martingale difference process where
 493 $\mathbb{E}(\Xi_n | \mathcal{F}_{t_n}) = 0$, and for integers $p \geq 2$ and a constant $s \in \mathbb{R}$ there are finite and
 494 strictly positive constants c_1 and c_2 such that $\mathbb{E}(|\Xi_n|^p) \leq c_1\delta^{p(s+1/2)}$. The process
 495 Θ similarly has $\mathbb{E}(|\Theta_n|^p) \leq c_2\delta^{p(s+1)}$. Then there exists constants c_3 and c_4 which
 496 depends only on L_A , L_B , and p such that $\mathbb{E}(\sup_{n \leq N} |\mathcal{E}_n|^p) \leq c_3(c_4c_1 + c_2)\delta^{ps}$, where
 497 $c_4 = 18p^{3/2}(p-1)^{-3/2}$.*

498 *Proof.* The proof is given by Giles and Sheridan-Methven [15, lemma 4.3] [49,
 499 lemma 5.2.3], and proceeds by a combination of Jensen's inequality, the discrete
 500 Burkholder-Davis-Gundy inequality, and the discrete Grönwall inequalities. \square

501 By considering the difference between the process \bar{X}_n calculated in high and
 502 low precision, then the result $\mathbb{E}(|\hat{X}_N - \bar{X}_N|^2) \leq CN\varrho^2$ from Arciniega and Allen
 503 [1, theorem 2.2] immediately follows from Lemma 2.5. Furthermore, we obtain an
 504 identical bound from Lemma 2.5 for the nett error that arises from Model 2.1.

505 LEMMA 2.6. *Using Model 2.1 for the rounding errors, then $\mathbb{E}(|\hat{X}_N - \bar{X}_N|^2) \leq$
 506 $CN\varrho^2$.*

507 *Proof.* Defining $\mathcal{E}_n := \hat{X}_n - \bar{X}_n$ and differencing the appropriate Euler-Maruyama

508 schemes for \widehat{X}_n and \overline{X}_n we obtain

$$\begin{aligned}
509 \quad \mathcal{E}_{n+1} &= \mathcal{E}_n + \delta \underbrace{(a(t_n, \widehat{X}_n) - a(t_n, \overline{X}_n))}_{\mathcal{A}_n} + \sqrt{\delta} \widetilde{Z}_n \underbrace{(b(t_n, \widehat{X}_n) - b(t_n, \overline{X}_n))}_{\mathcal{B}_n} \\
510 \quad &+ \underbrace{\sqrt{\delta} b(t_n, \widehat{X}_n) (\widehat{Z}_n - \widetilde{Z}_n) - \eta_n}_{\Xi_n} - \underbrace{\eta'_n}_{\Theta_n}, \\
511
\end{aligned}$$

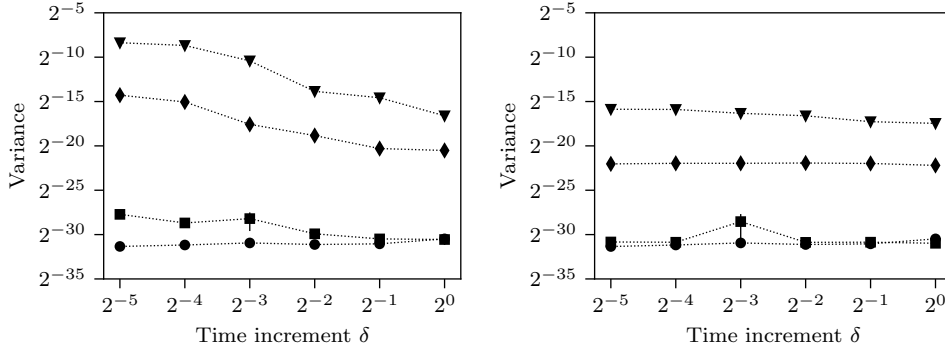
512 where we have indicated the equivalent terms in Lemma 2.5. The bounds on \mathcal{A}_n
513 and \mathcal{B}_n follow from the standard assumptions of a and b being spatially Lipschitz
514 continuous. Furthermore, for the Ξ_n term this is zero mean and has $\mathbb{E}(|\Xi_n|^p) \leq$
515 $O(\delta^{p/2}) + O(\varrho^p) \leq O(\varrho^p)$, corresponding to $s = -\frac{1}{2}$ in Lemma 2.5 and $c_1 \propto \varrho^p$.
516 Similarly, from Model 2.1 we have $\mathbb{E}(|\Theta_n|^p) \leq C \varrho^p \delta^{p/2}$, also corresponding to $s = -\frac{1}{2}$
517 and $c_2 \propto \varrho^p$. Thus from Lemma 2.5 we obtain $\mathbb{E}(|\widehat{X}_N - \overline{X}_N|^p) \leq O(\varrho^p \delta^{-p/2})$, which
518 when we set $p = 2$ and note that $\delta \propto \frac{1}{N}$ obtains the desired bound. \square

519 The significant insight provided by Lemma 2.6, which we saw in its proof when
520 we applied Lemma 2.5, is that the nett contributions from the η and η' processes grow
521 at the same rate. Although the η' process may be smaller than η by a factor of $\sqrt{\delta}$ in
522 Model 2.1, because it is not zero mean, its contributions do not cancel, and thus can
523 build up at a faster rate. Furthermore, this demonstrates that the updating process
524 is permitted a systematic rounding error process, provided it is $O(\sqrt{\delta})$.

525 The variance predicted by Lemma 2.6 is shown in Figure 2.3(a). For the approx-
526 imate random variables we have used the high fidelity piecewise cubic approximation
527 from Giles and Sheridan-Methven [16]. This is to ensure that the η_n term within
528 the Ξ_n process in the proof of Lemma 2.6 is the dominant term for moderately low
529 precisions. Using the mpmath Python library [28] we adjust the number of bits used
530 in the mantissa for the low precision approximate random variables. We use 7, 10,
531 and 23 bits, corresponding to the precisions for `bf16`, half, and single precisions
532 respectively. We also consider an artificial precision using 16 bits for the mantissa to
533 represent an intermediate precision level between half and single precision. For the
534 stochastic process we simulate a geometric Brownian motion where $a(t, X_t) \equiv \mu X_t$
535 and $b(t, X_t) \equiv \sigma X_t$ for strictly positive constant μ and σ . Following the setup from
536 Giles [13] we choose $\mu = 0.05$, $\sigma = 0.2$, and $X_0 = 1$.

537 We can see from Figure 2.3(a) that we approximately observe the growth in
538 the variance for half precision anticipated by Lemma 2.6. However, the brain float
539 precision seems to exhibit a variance closer to the worst case bound from Omland
540 [42]. Interestingly, we see that the higher precision intermediate and single precision
541 results appear to exhibit an approximate $O(1)$ variance. The reason for this is because
542 at such high precisions, the η term within Ξ_n is the smaller of the two terms, with
543 the approximation error $Z_n - \widetilde{Z}_n$ being the more dominant error. As shown by
544 Giles and Sheridan-Methven [15] this produces a resultant error independent of the
545 discretisation δ , and so appears $O(1)$, as observed. Furthermore, the intermediate
546 and single precision variances are approximately the same values, again indicating
547 that this error is dominated by the approximate random variables' fidelity rather
548 than their floating point precision.

549 **2.3. Kahan compensated summation.** The Euler-Maruyama scheme con-
550 sists of performing a cumulative summation of a sequence of update terms. The
551 problem of summing a sequence of floating point numbers and minimising the cu-
552 mulative rounding error is well known, and there have been a variety of methods
553 developed to overcome this, such as *pair wise summation* or *compensated summation*



(a) Without Kahan compensated summation.

(b) With Kahan compensation summation.

FIG. 2.3. The variance of the difference $\widehat{X}_T - \overline{X}_T$ between the exact Euler-Maruyama estimate \widehat{X} and the estimate using low precision approximate random variables \overline{X} from a piecewise cubic approximation with different precisions, corresponding to Lemma 2.6. The precisions use differing numbers of bits for the mantissa: (\blacktriangledown) 7, (\blacklozenge) 10, (\blacksquare) 16, and (\bullet) 23 bits. Results are shown with and without Kahan compensated summation, which is presented in Section 2.3.

554 [22, 4.1]. As the Euler-Maruyama scheme is a sequential and incremental procedure, a
555 compensated summation is an appropriate technique. As our schemes are motivated
556 by computational speed, then *Kahan compensated summation* [29], being the least
557 numerically intensive, is the most suitable candidate for incorporating into the Euler-
558 Maruyama scheme. The Kahan compensated summation adds a given increment, and
559 then subtracts away the computed summation prior to that increment. The differ-
560 ence of this inferred increment from the original provides an estimate for the incurred
561 rounding error, which is then adjusted for when adding subsequent terms in the se-
562 quence. The Kahan compensated summation procedure is outlined in Algorithm 2.1,
563 and a C implementation appropriate for incorporating this into the Euler-Maruyama
564 scheme in single precision is shown in Code 1. Interestingly, the related use of Kahan
565 compensated summation in the numerical solution of ordinary differential equations
566 was first proposed by Vitasek [52], and is demonstrated by Higham [21, pages 86–87].

Algorithm 2.1 Kahan compensated summation.

Input: A sequence $\{x_1, x_2, \dots, x_N\}$ of N floating point numbers.

Output: A high accuracy estimate of the summation $\sum_{i=1}^N x_i$.

- 1: Initialise both an accumulator a and compensation c to zero.
 - 2: **for all** $x_i \in \{x_1, x_2, \dots, x_n\}$ **do**
 - 3: Calculate a compensated increment $y \leftarrow x_i - c$.
 - 4: Keep the original $a_{\text{original}} \leftarrow a$.
 - 5: Add the compensated increment $a_{\text{new}} \leftarrow a + y$.
 - 6: Update the compensation $c \leftarrow (a_{\text{new}} - a_{\text{original}}) - y$.
 - 7: Update the accumulator $a \leftarrow a_{\text{new}}$.
 - 8: **end for**
 - 9: Use a to estimate $\sum_{i=1}^N x_i$.
-

567 An error analysis of Kahan compensated summation is provided by Higham [21,
568 page 791, (3.11)], Knuth [33, Exercise 19, pages 229 and 571–573], and Goldberg [19].

```

float compensated_EM_scheme(float X, float dX, float * comp)
{
    float compensated_increment = dX - (*comp);
    float accumulated_sum = X + compensated_increment;
    (*comp) = (accumulated_sum - X) - compensated_increment;
    return accumulated_sum;
}

```

CODE 1

C implementation of the Euler-Maruyama scheme using the Kahan compensated summation, implementing lines 3–7 from Algorithm 2.1.

569 The Kahan compensated summation shown in Algorithm 2.1 has the overall absolute
570 and relative error bounds of

$$571 \quad (2\varrho + O(N\varrho^2)) \sum_{i=1}^N |x_i| \quad \text{and} \quad (2\varrho + O(N\varrho^2)) \frac{\sum_{i=1}^N |x_i|}{\left| \sum_{i=1}^N x_i \right|}$$

572 respectively. The ratio of $\sum_{i=1}^N |x_i|$ to $|\sum_{i=1}^N x_i|$ is known as the condition number,
573 representing the sensitivity of the summation to rounding error [22, 50]. For a series
574 of zero mean random variables, the condition number can be expected to be $O(\sqrt{N})$,
575 but for several stochastic process (e.g. geometric Brownian motion) the drift term
576 causes the increments to have a non zero mean, and hence the condition number
577 approximately limits to a constant.

578 Based on the usual error analysis of Kahan compensated summation, one might
579 expect that the leading order error from a Kahan compensated Euler-Maruyama
580 scheme should mostly have an $O(1)$ error dependence on N , until eventually the
581 higher order term takes effect for sufficiently large N . However, inspecting Model 2.1
582 and Code 1 we see that the Kahan compensated summation is designed to tackle
583 the leading order η term. However, in computing the Euler-Maruyama update, the
584 smaller η' error is not compensated for, and thus will persist. Thus, even with Kahan
585 compensated summation, we see from Lemma 2.6 that we still expect a nett leading
586 order rounding error $O(\sqrt{N}\varrho)$. Overall then, we see that as we increase N , we ex-
587 pect for small N an $O(\varrho)$ error, for very large N an $O(N\varrho^2)$ error, and possibly an
588 intermediate $O(\sqrt{N}\varrho)$ error.

589 The variances from the approximations obtained by incorporating a Kahan com-
590 pensated summation into the modified Euler-Maruyama scheme are shown in Fig-
591 ure 2.3(b). The first thing to notice from this is that all the variances appear to be
592 $O(1)$, in keeping with the leading order error anticipated. For the brain float and
593 half precision variances, there is a separation between these which is approximately a
594 factor of 2^{-6} . As we expect the leading order error to be $O(\varrho)$, then we expect a re-
595 duction in the variance of approximately $\frac{\varrho_{10}^2}{\varrho^2} \approx 2^{-6}$, where we have used the subscript
596 to denote the number of mantissa bits. Hence we can see the reduction in variance
597 is approximately as anticipated. As for the intermediate and single precisions, we see
598 these cluster on top of each other, again indicating that the dominant error is from
599 approximating the Gaussian distribution rather than from finite precision arithmetic.

600 Our incorporation of the Kahan compensated summation scheme has been qual-
601 ified by heuristics and our empirical results. An exact and detailed analysis of the
602 errors remaining after incorporating Kahan compensated summation with Model 2.1

603 remains open for future analysis. The key point we highlight is its high level utility
604 and interpretation, which we further utilise and explore when discussing multilevel
605 Monte Carlo settings in Section 3, where it will considerably improve the capabilities
606 of working with half precision floating point numbers.

607 **3. Multilevel Monte Carlo.** We now have two possible types of simulation: an
608 expensive but precise simulation using exact Gaussian random variables in a high float-
609 ing point precision, and a cruder and cheaper simulation using approximate Gaussian
610 random variables in low precision. The accuracy of the former can be combined with
611 the speed of the latter by a multilevel Monte Carlo formulation [13, 14].

612 As a brief review of multilevel Monte Carlo, and how our work fits within this,
613 let us suppose we wish to compute the expectation of some functional P which acts
614 on the terminal solution X_T of the underlying stochastic process. The simulations
615 can be performed using various levels of temporal discretisation, where we index the
616 levels by l . We suppose there are $L + 1$ levels such that $l \in \{0, 1, 2, \dots, L\}$, where
617 $l = 0$ corresponds to the coarsest possible discretisation, and $l = L$ the finest. The
618 approximation coming from the usual Euler-Maruyama scheme for a particular level
619 l we denote by \hat{P}_l . Additionally, we denote those arising from the modified Euler-
620 Maruyama scheme using high precision approximate random variables by \tilde{P}_l , and low
621 precision approximate random variables by \bar{P}_l . For notational simplicity we use the
622 convention $\hat{P}_{-1} := \tilde{P}_{-1} := \bar{P}_{-1} := 0$. Giles and Sheridan-Methven [15, 16] suggest
623 incorporating the approximate random variables using the nested multilevel Monte
624 Carlo framework

$$\begin{aligned}
625 \quad \mathbb{E}(P) &\approx \mathbb{E}(\hat{P}_L) = \sum_{l=0}^L \mathbb{E}(\hat{P}_l - \hat{P}_{l-1}) \\
626 &= \sum_{l=0}^L \mathbb{E}(\tilde{P}_l - \tilde{P}_{l-1}) + \mathbb{E}(\hat{P}_l - \hat{P}_{l-1} - \tilde{P}_l + \tilde{P}_{l-1}) \\
627 &= \sum_{l=0}^L \mathbb{E}(\bar{P}_l - \bar{P}_{l-1}) + \mathbb{E}(\hat{P}_l - \hat{P}_{l-1} - \bar{P}_l + \bar{P}_{l-1}), \\
628
\end{aligned}$$

629 where the first approximation is the regular Monte Carlo procedure [18], the first
630 equality is the usual multilevel Monte Carlo decomposition [13], the second equality
631 is the nested multilevel Monte Carlo framework [15, 16], and the final equality is
632 the same nested multilevel framework utilising low precision approximate random
633 variables.

634 Importantly, for a given level l , the fine path's discretisation δ^f and the coarse
635 path's δ^c are given by $\delta^f = 2^{-l}$ and $\delta^c = 2\delta^f$ respectively. The coarse path's Wiener
636 increments are produced by the pairwise summation of the fine path's Wiener in-
637 crements. Furthermore, the Wiener increments ΔW are produced using Gaussian
638 increments Z where $\Delta W := \sqrt{\delta}Z$, and the Gaussian increments are produced using
639 the inverse transform method where $Z_n = \Phi^{-1}(U_n)$. Crucially, the exact random vari-
640 ables Z_n and approximate random variables \tilde{Z}_n are tightly coupled by ensuring they
641 are both generated using the same underlying random variables, where $Z_n = \Phi^{-1}(U_n)$
642 and $\tilde{Z}_n = \tilde{\Phi}^{-1}(U_n)$, with U_n being the same for both.

643 Giles and Sheridan-Methven [16] consider the usual multilevel estimator $\hat{\theta}$ and

644 the nested multilevel estimator $\bar{\theta}$ where

$$645 \quad \hat{\theta} := \sum_{l=0}^L \frac{1}{\hat{m}_l} \sum^{\hat{m}_l} (\hat{P}_l - \hat{P}_{l-1})$$

646 and

$$648 \quad \bar{\theta} := \sum_{l=0}^L \left(\frac{1}{\bar{m}_l} \sum^{\bar{m}_l} (\bar{P}_l - \bar{P}_{l-1}) + \frac{1}{\bar{M}_l} \sum^{\bar{M}_l} (\hat{P}_l - \hat{P}_{l-1} - \bar{P}_l + \bar{P}_{l-1}) \right),$$

649 where \hat{m}_l , \bar{m}_l , and \bar{M}_l are the number of paths generated, each with a computational
650 cost of \hat{c}_l , \bar{c}_l , and \bar{C}_l , and variance \hat{v}_l , \bar{v}_l , and \bar{V}_l respectively. Letting \hat{T} denote the
651 total computational time to achieve a mean squared error ε^2 using the estimator $\hat{\theta}$,
652 and similarly \bar{T} using $\bar{\theta}$, then Giles and Sheridan-Methven [16] show

$$654 \quad \hat{T} = 2\varepsilon^{-2} \left(\sum_{l=0}^L \sqrt{\hat{v}_l \hat{c}_l} \right)^2 \quad \text{and} \quad \bar{T} = 2\varepsilon^{-2} \left(\sum_{l=0}^L \sqrt{\bar{v}_l \bar{c}_l} + \sqrt{\bar{V}_l \bar{C}_l} \right)^2,$$

655 and hence an overall temporal saving of

$$656 \quad \bar{T} \approx 2\varepsilon^{-2} \left(\sum_{l=0}^L \sqrt{\hat{v}_l \hat{c}_l} \left(\sqrt{\frac{\bar{v}_l \bar{c}_l}{\hat{v}_l \hat{c}_l}} + \sqrt{\frac{\bar{V}_l \bar{C}_l}{\hat{v}_l \hat{c}_l}} \right) \right)^2 \leq \hat{T} \max_{l \leq L} \left\{ \frac{\bar{v}_l \bar{c}_l}{\hat{v}_l \hat{c}_l} \left(1 + \sqrt{\frac{\bar{V}_l \bar{C}_l}{\bar{v}_l \bar{c}_l}} \right)^2 \right\}.$$

657 When the approximation's fidelity is such that $\bar{v}_l \approx \hat{v}_l$, the term $\frac{\bar{v}_l \bar{c}_l}{\hat{v}_l \hat{c}_l} \approx \frac{\bar{c}_l}{\hat{c}_l}$ measures
658 the potential time savings, and the term $(1 + (\bar{V}_l \bar{C}_l / \bar{v}_l \bar{c}_l)^{1/2})^2$ assesses the efficiency
659 at realising these savings. Achieving a balance between these two is required, where
660 the approximation should be sufficiently fast so there is the potential for large savings,
661 but of a sufficient fidelity so the variance of the more expensive four way difference
662 $\hat{P}_l - \hat{P}_{l-1} - \bar{P}_l + \bar{P}_{l-1}$ is considerably lower than the variance of the cheaper two way
663 difference $\hat{P}_l - \hat{P}_{l-1}$.

664 Giles and Sheridan-Methven [15, 16] investigate the variance of this final four
665 way difference when using approximate random variables assuming infinite precision
666 arithmetic for a variety of functional types [15, lemmas 4.10 and 4.11] [49, corollar-
667 ies 6.2.6.2 and 6.2.6.3]. They find that for Lipschitz continuous and differentiable
668 functionals that

$$669 \quad (3.1) \quad \|\hat{P}_l - \hat{P}_{l-1} - \tilde{P}_l + \tilde{P}_{l-1}\|_p \leq O(\delta^{1/2} \|Z - \tilde{Z}\|_{p'})$$

670 for some p' such that $2 \leq p < p' < \infty$. For simplicity we will restrict our attention
671 to functionals which are also Lipschitz continuous and differentiable, and we consider
672 just the underlying process and take $P(X) \equiv X$.

673 In an analogous manner to Lemma 2.6 we can consider the variance of this four
674 way difference arising in the nested multilevel Monte Carlo framework when the pre-
675 cision is simultaneously lowered when switching to approximate random variables,
676 giving rise to Lemma 3.1.

677 **LEMMA 3.1.** *For fine and coarse path simulations constructed using approximate*
678 *random variables as described by Giles and Sheridan-Methven [15], then with rounding*
679 *errors described by Model 2.1 we have*

$$680 \quad \mathbb{E}(|\hat{X}_l - \hat{X}_{l-1} - \bar{X}_l + \bar{X}_{l-1}|^2) \approx O(\delta \mathbb{E}(|Z - \tilde{Z}|^{2+\epsilon})) + O(\delta^{-1} \rho^2)$$

681 as the discretisation δ parameter decreases for some $\epsilon \in (0, \infty)$.

682 *Proof.* The initial $O(\delta \mathbb{E}(|Z - \tilde{Z}|^{2+\epsilon}))$ term comes from (3.1) [15, 16] taking
 683 $P(X) \equiv X$ and $p = 2$. The second $O(\delta^{-1} \varrho^2)$ term comes from the η and η' con-
 684 tributions from Model 2.1, and their final nett contributions arise using Lemma 2.5
 685 in an identical manner to the proof of Lemma 2.6. \square

686 We can compare the realised variances predicted by Giles and Sheridan-Methven
 687 [15, 16] when using approximate random variables in infinite precision, with those
 688 predicted by Lemma 3.1. We use the piecewise linear approximation by Giles and
 689 Sheridan-Methven [16], and use half precision capable hardware (rather than emula-
 690 tion using the mpmath Python library). We use C code on a Nvidia Jetson AGX
 691 Xavier machine, containing a Nvidia 12 core Volta GPU and an 8 core Arm v8.2 64
 692 bit CPU. Both the CPU and the GPU support half precision floating point arithmetic
 693 in hardware, and so we run the code on the CPU using the `_Float16` data type for half
 694 precision, compiled with `gcc` and notably with the flags `-O0` and `-march=armv8.2-a+fp16`.
 695 The first flag ensures no compiler optimisations are issued, guaranteeing the Kahan
 696 compensated summation is not removed by the compiler, and the second ensures
 697 half precision data types and operations are accessible and used. The results for the
 698 variances of the various multilevel Monte Carlo terms are shown in Figure 3.1.

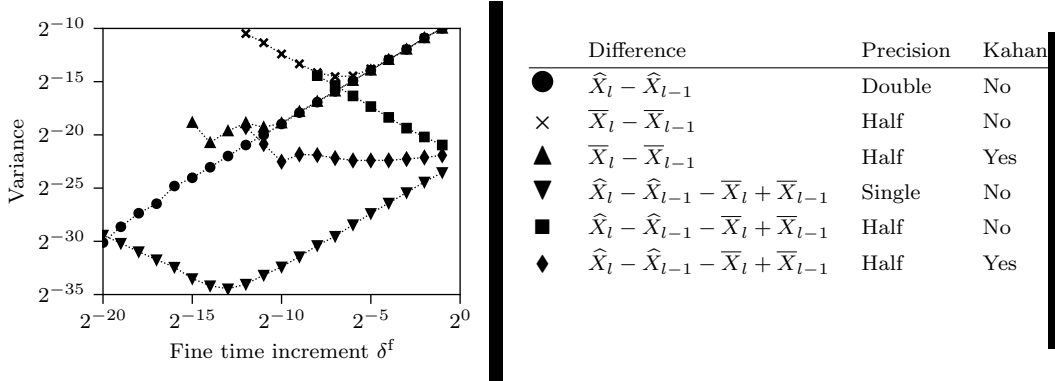


FIG. 3.1. The variance reduction when switching to approximate random variables with a nested multilevel Monte Carlo framework, showing the variances of various two and four way differences. The key indicates the difference, precision, and whether Kahan compensated summation is used.

699 Inspecting Figure 3.1 we can make several observations. The first is that the usual
 700 two way difference (computed in double precision) exhibits the usual $O(\delta)$ variance
 701 decay rate, as is to be expected and is a standard result [18, 32].

702 The next item of interest is the behaviour of the four way difference computed in
 703 single precision. Down to discretisations as fine as $\delta \approx 2^{-13}$ we can see the variance
 704 decays at the same $O(\delta)$ rate, as already predicted and demonstrated by Giles and
 705 Sheridan-Methven [16, 4.1]. However, the novel feature predicted by Lemma 3.1
 706 is the emergence of the $O(\delta^{-1})$ rate for very fine discretisation, arising from Model 2.1.
 707 We can see that the onset of rounding error is not immediately catastrophic, and that
 708 for $\delta \approx 2^{-17}$ there is still a reduction in the variance between the two and four way
 709 differences which is still approximately 2^{-6} . However, eventually, for discretisations as
 710 fine as $\delta \approx 2^{-20}$ there is no reduction in variance. Ultimately, this demonstrates that
 711 double precision is typically superfluous for Monte Carlo path simulations (ignoring

712 sensitivity calculations for computing derivatives by finite differences), and single
 713 precision is sufficient.

714 The main items of particular interest are the half precision variances. For sim-
 715 plicity we begin by inspecting the difference without Kahan compensated summation.
 716 The variance of the two way difference term decreases in line with the double preci-
 717 sion two way difference down to approximately $\delta \approx 2^{-7}$, and thereafter the effects of
 718 rounding error become dominant. This immediately places a lower limit on an un-
 719 compensated half precision framework, which is valid for discretisations coarser than
 720 $\delta \geq 2^{-7}$. Considering the four way difference, we see that at the very coarsest level
 721 there is an appreciable variance reduction by a factor of approximately 2^{-12} , although
 722 rounding error has already started to become dominant. We see that as the discreti-
 723 sations become finer the rounding error increases, and at $\delta \approx 2^{-7}$ coincides with the
 724 two way difference. However, for discretisations coarser than $\delta \geq 2^{-4}$, there is still at
 725 least a variance drop by approximately a factor of 2^{-6} . This suggests that while half
 726 precision calculations may be fast, if the rounding error is not compensated for, then
 727 they are only useful on the coarsest few levels.

728 If we incorporate a Kahan compensated summation to the Euler-Maruyama scheme
 729 when using half precision approximate random variables, the picture improves. The
 730 first item to note is the variance of the two way difference, which mirrors the double
 731 precision’s two way variance down to approximately $\delta \approx 2^{-11}$, placing a lower limit
 732 on the minimum possible discretisation with $\delta \geq 2^{-11}$. For the four way difference,
 733 at the very coarsest level we see approximately the same variance as the uncompen-
 734 sated half precision four way difference, as might be expected. However, with the
 735 compensation, the error is an $O(1)$ constant as the discretisation becomes ever finer
 736 down to $\delta \approx 2^{-10}$. For discretisations $\delta < 2^{-10}$ a higher order error process appears
 737 to dominate, and for such fine discretisations the two way and four way variances
 738 coincide. This suggests that half precision simulations using Kahan compensated
 739 summation are applicable for much finer discretisation than equivalent simulations
 740 without the Kahan compensated summation. A variance reduction of approximately
 741 2^{-6} is achieved for discretisations $\delta \geq 2^{-8}$. Thus Kahan compensated summation ap-
 742 proximately doubles the scope of practical applicability for half precision simulations.

TABLE 3.1

*Performance of various exact and approximate implementations of the inverse Gaussian cu-
 mulative distribution function, and the possible speed ups offered.*

(a) The time to generate exact and approximate Gaussian random variables.				(b) The maximum possible speed ups.		
Description	Precision	Clock cycles	Source	Precision	Kahan	Speed up
Intel (HA)	Single	3.5 ± 0.2	[16]	Single	No	7
Piecewise linear	Single	0.5 ± 0.1	[16]	Half	No	14
Piecewise linear	Half	0.25	Speculated	Half	Yes	10

743 We take the timing results on Intel AVX-512 Skylake hardware from Giles and
 744 Sheridan-Methven [16]. For the exact Gaussian distribution, we take as our baseline
 745 the single precision Intel high accuracy (HA) function. In lieu of vectorised half
 746 precision capable hardware, we speculate that for the approximate random variables,
 747 that half precision input can be processed in half the time as single precision input.
 748 Overall then we have the times shown in Table 3.1(a). Furthermore, we make the
 749 idealised assumption that the cost of the simulations is entirely based on the cost of
 750 generating random numbers, and neglect the cost of all other arithmetic operations.

751 These give the maximum potential speed ups shown in Table 3.1(b). For the half
 752 precision approximation using Kahan compensated summation, we again speculate
 753 and suggest an intermediate value between the single and half precision offerings.

754 For each possible discretisation level, the estimated speed ups predicted for each
 755 individual level from the nested multilevel Monte Carlo analysis are shown in Fig-
 756 ure 3.2. We can see from this that for the bulk of levels the single precision approxi-
 757 mation offers a good potential speed up by a factor of approximately 7. However, for
 758 the very coarsest few levels, the half precision approximations without Kahan com-
 759 pensated summation offer superior speed ups by a factor of 10–12. Dependent on the
 760 speed reduction that comes from incorporating the half precision Kahan compensated
 761 summation, there is the possibility of these offering a third intermediate regime where
 762 they are the optimal choice. For our speculated speed ups from Table 3.1(b) we see
 763 there appears to be such an intermediate region, although it is not overwhelmingly
 764 competitive compared to the two other alternatives.

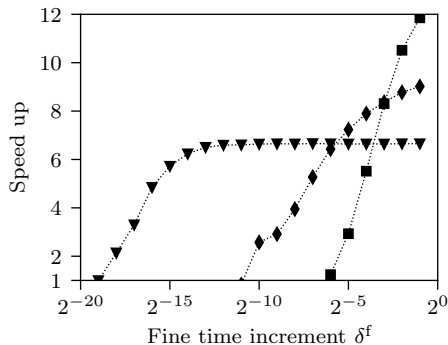


FIG. 3.2. The potential savings from a nested multilevel Monte Carlo framework using approximate random variables from a piecewise linear approximation for various discretisation levels. (▼) Single precision. (■) Half precision without Kahan compensated summation. (◆) Half precision with Kahan compensated summation.

765 It is worth remarking that when using the Euler-Maruyama scheme, as the scheme
 766 has a strong convergence order of $\frac{1}{2}$, the computational work load is approximately
 767 spread evenly over the various levels [13]. However, for numerical schemes with higher
 768 strong convergence orders, such as the Milstein scheme which has order 1 strong
 769 convergence [18, 32], the work load is predominantly concentrated on the coarsest
 770 levels [13]. The implication of this is that the potential half precision speed ups
 771 offered on the coarsest levels, even without Kahan compensated summation, may
 772 well dominate the multilevel savings. Hence, while half precision appears attractive
 773 even with multilevel Monte Carlo frameworks using the Euler-Maruyama scheme, this
 774 becomes even more so for the Milstein scheme and other higher order methods.

775 Currently, the use of approximate random variables and their incorporation in
 776 stochastic simulations and nested multilevel Monte Carlo frameworks only has rigor-
 777 ous supporting analytic results for the Euler-Maruyama scheme [15, 16]. Extensions
 778 to higher order schemes, such as the Milstein scheme, have only been discussed and
 779 demonstrated empirically [15, 16]. Consequently, there is currently insufficient sup-
 780 porting analysis to attempt to marry low precision approximate random variables in
 781 a nested multilevel Monte Carlo framework with higher order numerical schemes at
 782 an analytic level, and such analysis remains open for further research.

783 Lastly, we can speculate about the utility of brain floats compared to regular half
784 precision floats (taken to be the IEEE specification). Brain floats and half precision
785 floats are both 16 bits in size, but differ in their trade off between precision and range.
786 Brain floats have a much larger range and lower precision than regular half precision.
787 Having a lower precision will likely mean that the initial impact of rounding error will
788 be even more severe than for half precision. This means the variance reduction will
789 be less favourable, and the nested multilevel Monte Carlo framework less efficient on
790 each level, and useful over even fewer discretisations. Without cause to suspect brain
791 floats will be any faster than regular half precision, this would suggest that while half
792 precision is attractive for multilevel Monte Carlo applications, and brain floats may
793 similarly be attractive for the same reasons, there is no reason to speculate that brain
794 floats will be competitive over regular half precision floats.

795 **4. Conclusions.** Performing calculations in high precisions may assuage worries
796 about rounding errors, but makes several computations needlessly expensive. Consid-
797 ering the numerical simulation of stochastic differential equations, based on previous
798 work using computationally cheap approximate random variables, we couple their
799 incorporation into the Euler-Maruyama scheme with low precision implementations.
800 We develop a new model for the nett rounding error incurred which allows for both
801 systematic and unsystematic errors, analysing how the two can be balanced under ap-
802 propriate assumptions. Kahan compensated summation is also discussed as a means
803 of removing the leading order rounding error. This rounding error model is incor-
804 porated into a nested multilevel Monte Carlo scheme allowing for the speed of low
805 precisions to be capitalised on without losing accuracy, finding that single precision
806 is applicable for most discretisations and offers good potential speed ups, while half
807 precision appears to offer superior speeds at only the coarsest few discretisation levels.

808 Introducing finite precision calculations, we discuss the appeal of working in ever
809 lower precisions, and the attraction of using half precision in various applications.
810 Low precision offers improved speed by increasing bandwidth and decreasing floating
811 point calculation times, and is rapidly gaining traction in software and hardware,
812 primarily due to applications in machine learning. However, low precision comes
813 with appreciable finite precision rounding error, whose effects are felt in applications
814 including: linear algebra, machine learning, stochastic simulation, and various others.
815 To mitigate against this there are high precision libraries, compensated summation
816 schemes, and mathematical techniques such as Richardson extrapolation, although
817 the problem of rounding error can be severe for 16 bit floating point formats such as
818 half precision or brain floats.

819 Considering the setting of stochastic simulations and approximating solutions
820 of stochastic differential equations using the Euler-Maruyama scheme we introduce
821 Model 2.1 as a novel description for the effects of rounding error. Similar previous
822 models for rounding error in this setting are notably from Arciniega and Allen [1]
823 and Omland [42], for the average and worst case scenarios respectively. The model
824 we introduce is based on the framework by Omland [42] and recovers a similar model
825 to the one by Arciniega and Allen [1]. However, comparing our model to that from
826 Arciniega and Allen [1], ours has several benefits. Firstly, ours facilitates the incor-
827 poration of approximate random variables, which Giles and Sheridan-Methven [16]
828 showed offers considerable potential speed improvements. Our model now rigorously
829 justifies a leading order zero mean contribution, and permits a smaller second order
830 and possibly non zero mean contribution, thus also facilitating systematic errors. Us-
831 ing Lemma 2.5 we show the nett contributions from these two terms grow at the same

832 rate for the Euler-Maruyama scheme, providing novel insight into the size permissible
833 for systematic errors.

834 Reviewing Kahan compensated summation, we discuss the leading order constant
835 error this can be expected to produce. We numerically simulated geometric Brownian
836 motion processes, finding the errors predicted by Lemma 2.6 are observed empirically,
837 alongside the anticipated error resulting from incorporating Kahan compensated sum-
838 mation. Furthermore, while Kahan compensated summation cancels the leading order
839 rounding error, there is a balance between this and the error resulting from using ap-
840 proximate random variables, which becomes the dominant source of error in higher
841 fidelity simulations.

842 Low precision implementations of approximate random variables are best utilised
843 with a nested multilevel Monte Carlo framework, paralleling the setup by Giles and
844 Sheridan-Methven [16]. Crucially, for practitioners working in high performance sci-
845 entific computing, this integration both opens up the capabilities of realising the
846 potential speed ups offered from low precision calculations and also gives analytic
847 assurances that accuracy need not be lost. The nested multilevel Monte Carlo setup
848 produces a four way difference, whose variance we predict in Lemma 3.1 and empiri-
849 cally observe. As expected, we find the errors resulting from introducing approximate
850 random variables and also from low precision calculations are orthogonal effects. Us-
851 ing the nested multilevel Monte Carlo framework, we calculate that for a very wide
852 range of possible discretisation levels that single precision simulations offer a speed
853 up by a factor of approximately 7, as already shown by Giles and Sheridan-Methven
854 [16]. However, for the very coarsest few levels, we demonstrated that half precision
855 calculations, despite incurring significant rounding error from the very onset, can be
856 successfully utilised within a multilevel Monte Carlo scheme. On these coarsest few
857 levels, half precision can offer speed improvements by a factor of 10–12. For the
858 Euler-Maruyama scheme the work load is usually evenly spread across the various
859 levels, whereas for the higher order Milstein scheme the work load is dominated by
860 the coarsest few levels. Thus we have been able to demonstrate the utility and ap-
861 plicability of half precision approximate random variables for stochastic simulation
862 applications. Lastly, dependent on the cost model for the added arithmetic required
863 for Kahan compensated summation, we were also able to demonstrate that this may
864 provide a further intermediate region, extending the range of half precision to even
865 finer discretisation levels, compounding their benefits.

866 All of the code to produce these figures is freely available and hosted by Sheridan-
867 Methven [48].

868 **5. Acknowledgements.** We would like to acknowledge and thank those who
869 have financially sponsored this work. This includes the Engineering and Physical Sci-
870 ences Research Council (EPSRC) and Oxford University’s centre for doctoral training
871 in Industrially Focused Mathematical Modelling (InFoMM), with the EP/L015803/1
872 funding grant. Furthermore, this research stems from a PhD project [49] which was
873 funded by Arm and NAG. Funding was also provided by the EPSRC ICONIC pro-
874 gramme grant EP/P020720/1, the Hong Kong Innovation and Technology Commis-
875 sion (InnoHK Project CIMDA), and by Mansfield College, Oxford.

876 **References.**

- 877 [1] Armando Arciniega and Edward Allen. Rounding error in numerical solution
878 of stochastic differential equations. *Stochastic analysis and applications*, 21(2):
879 281–300, 2003.

- 880 [2] Ivo Babuška. Numerical stability in mathematical analysis. In *IFIP congress (I)*,
881 volume 68, pages 11–23, 1968.
- 882 [3] Stephen J. Blundell and Katherine M. Blundell. *Concepts in thermal physics*.
883 Oxford University press, 2nd edition, 2014.
- 884 [4] George E.P. Box and Mervin E. Muller. A note on the generation of random
885 normal deviates. *The annals of mathematical statistics*, 29:610–611, 1958.
- 886 [5] Christian Brugger, Christian de Schryver, Norbert Wehn, Steffen Omland, Mario
887 Hefter, Klaus Ritter, Anton Kostiuk, and Ralf Korn. Mixed precision multi-
888 level Monte Carlo on hybrid computing systems. In *2104 IEEE conference on*
889 *computational intelligence for financial engineering & economics (CIFER)*, pages
890 215–222. IEEE, 2014.
- 891 [6] Neil Burgess, Chris Goodyer, Chris Hinds, and David Lutz. High-precision an-
892 chored accumulators for reproducible floating-point summation. *IEEE transac-*
893 *tions on computers*, 2018.
- 894 [7] Neil Burgess, Jelena Milanovic, Nigel Stephens, Konstantinos Monachopoulos,
895 and David Mansell. Bfloat16 processing for neural networks. In *2019 IEEE 26th*
896 *symposium on computer arithmetic (ARITH)*, pages 88–91. IEEE, 2019.
- 897 [8] Gary Chun Tak Chow, Anson Hong Tak Tse, Qiwei Jin, Wayne Luk, Philip H.W.
898 Leong, and David B. Thomas. A mixed precision Monte Carlo methodology for
899 reconfigurable accelerator systems. In *Proceedings of the ACM/SIGDA interna-*
900 *tional symposium on field programmable gate arrays*, pages 57–66, 2012.
- 901 [9] Matteo Croci and Michael B. Giles. Effects of round-to-nearest and stochastic
902 rounding in the numerical solution of the heat equation in low precision. *IMA*
903 *Journal of Numerical Analysis*, 43(3):1358–1390, 04 2022. ISSN 0272-4979. doi:
904 10.1093/imanum/drac012.
- 905 [10] Theodorus Jozef Dekker. A floating-point technique for extending the available
906 precision. *Numerische Mathematik*, 18(3):224–242, 1971.
- 907 [11] Sylvain Delattre and Jean Jacod. A central limit theorem for normalized func-
908 tions of the increments of a diffusion process, in the presence of round-off errors.
909 *Bernoulli*, 3(1):1–28, 1997.
- 910 [12] Laurent Fousse, Guillaume Hanrot, Vincent Lefèvre, Patrick Pélissier, and Paul
911 Zimmermann. MPFR: a multiple-precision binary floating-point library with
912 correct rounding. *ACM transactions on mathematical software (TOMS)*, 33(2):
913 13–es, June 2007.
- 914 [13] Michael B. Giles. Multilevel Monte Carlo path simulation. *Operations research*,
915 56(3):607–617, 2008.
- 916 [14] Michael B. Giles. Multilevel Monte Carlo methods. *Acta Numerica*, 24:259–328,
917 2015.
- 918 [15] Michael B. Giles and Oliver Sheridan-Methven. Analysis of nested multilevel
919 Monte Carlo using approximate normal random variables. *SIAM/ASA Journal*
920 *on Uncertainty Quantification*, 10(1):200–226, 2022.
- 921 [16] Michael B. Giles and Oliver Sheridan-Methven. Approximating inverse cumu-
922 lative distribution functions to produce approximate random variables. *ACM*
923 *transactions on mathematical software (TOMS)*, 49(3), September 2023.
- 924 [17] Michael B. Giles, Mario Hefter, Lukas Mayer, and Klaus Ritter. Random bit
925 multilevel algorithms for stochastic differential equations. *Journal of complexity*,
926 2019.
- 927 [18] Paul Glasserman. *Monte Carlo methods in financial engineering*, volume 53 of
928 *Stochastic modelling and applied probability*. Springer science & business media,
929 1st edition, 2013.

- 930 [19] David Goldberg. What every computer scientist should know about floating-point
931 arithmetic. *ACM computing surveys (CSUR)*, 23(1):5–48, 1991.
- 932 [20] Torbjörn Granlund and the GMP development team. GNU MP: the GNU mul-
933 tiple precision arithmetic library, 2012. URL <http://gmplib.org/>. Version 5.0.5.
- 934 [21] Nicholas J. Higham. The accuracy of floating point summation. *SIAM journal*
935 *on scientific computing*, 14(4):783–799, 1993.
- 936 [22] Nicholas J. Higham. *Accuracy and stability of numerical algorithms*, volume 80.
937 SIAM, 2nd edition, 2002.
- 938 [23] Nicholas J. Higham and Theo Mary. A new approach to probabilistic rounding
939 error analysis. *SIAM journal on scientific computing*, 41(5):A2815–A2835, 2019.
- 940 [24] Thomas E. Hull and J. R. Swenson. Tests of probabilistic models for propagation
941 of roundoff errors. *Communications of the ACM*, 9(2):108–113, February 1966.
942 ISSN 0001-0782.
- 943 [25] IEEE. IEEE standard for binary floating-point arithmetic, 1985. Computer
944 society standards committee. Working group of the microprocessor standards
945 subcommittee.
- 946 [26] IEEE. IEEE standard for binary floating-point arithmetic, 2008. Computer
947 society standards committee. Working group of the microprocessor standards
948 subcommittee.
- 949 [27] Ilse C.F. Ipsen and Hua Zhou. Probabilistic error analysis for inner products,
950 2019. arXiv:1906.10465.
- 951 [28] Fredrik Johansson et al. mpmath: a Python library for arbitrary-precision
952 floating-point arithmetic (version 0.18), December 2013. URL <http://mpmath.org/>.
- 953 [29] William Morton Kahan. Further remarks on reducing truncation errors. *Com-*
954 *munications of the association for computing machinery (ACM)*, 8:40, 1965.
- 955 [30] Dhiraj Kalamkar, Dheevatsa Mudigere, Naveen Mellempudi, Dipankar Das, Ku-
956 nal Banerjee, Sasikanth Avancha, Dharma Teja Vooturi, Nataraj Jammala-
957 madaka, Jianyu Huang, Hector Yuen, Jiyan Yang, Jongsoo Park, Alexander
958 Heinecke, Evangelos Georganas, Sudarshan Srinivasan, Abhisek Kundu, Misha
959 Smelyanskiy, Bharat Kaul, and Pradeep Dubey. A study of BFLOAT16 for deep
960 learning training, 2019. URL <https://arxiv.org/abs/1905.12322>.
- 961 [31] Andreas Klein. A generalized Kahan-Babuška summation algorithm. *Computing*,
962 76(3):279–293, 2006.
- 963 [32] Peter E. Kloeden and Eckhard Platen. *Numerical solution of stochastic differen-*
964 *tial equations*, volume 23 of *Stochastic modelling and applied probability*. Springer,
965 1999. Corrected 3rd printing.
- 966 [33] Donald E. Knuth. *The art of computer programming: seminumerical algorithms*,
967 volume 2. Addison-Wesley Professional, 2014.
- 968 [34] Peter Linz. Accurate floating-point summation. *Communications of the ACM*,
969 13(6):361–362, June 1970.
- 970 [35] Gurii Ivanovich Marchuk and Vladimir V. Shaidurov. *Difference methods and*
971 *their extrapolations*, volume 19. Springer science & business media, 2012.
- 972 [36] George Marsaglia and T.A. Bray. A convenient method for generating normal
973 variables. *SIAM Review*, 6(3):260–264, 1964.
- 974 [37] George Marsaglia and Wai Wan Tsang. The Ziggurat method for generating
975 random variables. *Journal of statistical software*, 5(1):1–7, 2000.
- 976 [38] Ole Møller. Quasi double-precision in floating point addition. *BIT numerical*
977 *mathematics*, 5(1):37–50, 1965.
- 978 [39] Mervin E. Muller. An inverse method for the generation of random normal devi-
979

- 980 ates on large-scale computers. *Mathematical tables and other aids to computation*,
981 12(63):167–174, 1958.
- 982 [40] Arnold Neumaier. Rounding error analysis of some methods for summing finite
983 sums (rundungsfehleranalyse einiger verfahren zur summation endlicher sum-
984 men). *ZAMM - Journal of Applied Mathematics and Mechanics (zeitschrift für*
985 *angewandte mathematik und mechanik)*, 54(1):39–51, 1974.
- 986 [41] Takeshi Ogita, Siegfried M Rump, and Shin’ichi Oishi. Accurate sum and dot
987 product. *SIAM journal on scientific computing*, 26(6):1955–1988, 2005.
- 988 [42] Steffen Omland. Mixed precision multilevel Monte Carlo algorithms for recon-
989 figurable computing systems, June 2016. PhD thesis/dissertation, D 386, Tech-
990 nische Universität Kaiserslautern (TUK).
- 991 [43] Steffen Omland, Mario Hefter, Klaus Ritter, Christian Brugger, Christian
992 de Schryver, Norbert Wehn, and Anton Kostiuik. Exploiting mixed-precision
993 arithmetics in a multilevel Monte Carlo approach on FPGAs. In *FPGA based*
994 *accelerators for financial applications*, pages 191–220. Springer, 2015.
- 995 [44] Lewis Fry Richardson. VIII. The deferred approach to the limit. *Philosophi-
996 cal transactions of the royal society of london. Series A, containing papers of a*
997 *mathematical or physical character*, 226(636-646):299–361, 1927.
- 998 [45] Mathieu Rosenbaum. Integrated volatility and round-off error. *Bernoulli*, 15(3):
999 687–720, 2009.
- 1000 [46] Oliver Sheridan-Methven. Approximating inverse cumulative distribution func-
1001 tions, 2020. URL [https://github.com/oliversheridanmethven/approximating-
1002 inverse_cumulative_distribution_functions](https://github.com/oliversheridanmethven/approximating-inverse_cumulative_distribution_functions). GitHub repository.
- 1003 [47] Oliver Sheridan-Methven. Getting started with approximate random vari-
1004 ables: a brief guide for practitioners, 2020. URL [https://github.com/
1005 oliversheridanmethven/approximate_random_variables](https://github.com/oliversheridanmethven/approximate_random_variables). GitHub repository.
- 1006 [48] Oliver Sheridan-Methven. Low precision approximate random vari-
1007 ables, 2020. URL [https://github.com/oliversheridanmethven/low_precision-
1008 approximate_random_variables](https://github.com/oliversheridanmethven/low_precision_approximate_random_variables). GitHub repository.
- 1009 [49] Oliver Sheridan-Methven. Nested multilevel Monte Carlo methods and a mod-
1010 ified Euler-Maruyama scheme utilising approximate Gaussian random variables
1011 suitable for vectorised hardware and low-precisions, 2021. DPhil. thesis, Mathe-
1012 matical Institute, University of Oxford.
- 1013 [50] Lloyd N. Trefethen and David Bau. *Numerical linear algebra*, volume 50. SIAM,
1014 3rd edition, 1997.
- 1015 [51] Warwick Tucker. *Validated numerics: a short introduction to rigorous computa-
1016 tions*. Princeton University press, 2011.
- 1017 [52] Emil Vitasek. The numerical stability in solution of differential equations. In
1018 *Conference on the numerical solution of differential equations*, pages 87–111.
1019 Springer, 1969.
- 1020 [53] James H. Wilkinson. Error analysis of direct methods of matrix inversion. *Journal*
1021 *of the ACM (JACM)*, 8(3):281–330, 1961.
- 1022 [54] James H. Wilkinson. Numerical linear algebra on digital computers. *IMA bulletin*,
1023 10(9/10):354–356, 1974.
- 1024 [55] James H. Wilkinson. Error analysis revisited. *IMA bulletin*, 22(11/12):192–200,
1025 1986.