

Stochastic Analysis and Correction of Floating Point Errors in Monte Carlo Simulations

Oliver Sheridan-Methven

27th February 2019

Supervisors:

Prof. Mike Giles	Oxford
Dr Christopher Goodyer	Arm



Engineering and Physical Sciences
Research Council

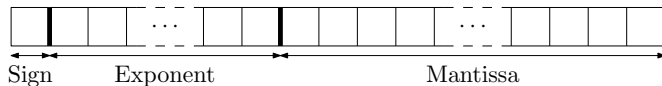
Overview



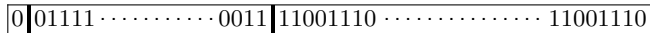
- 1 Fast computations and low-precision
- 2 Low-precision and higher speed
- 3 Are crude approximations useful?
- 4 Mixed-precision multilevel Monte Carlo
- 5 Convergence results
- 6 Recent results
- 7 Numerical results
- 8 Conclusions

$$x \equiv 0.0001101011010101 \dots$$

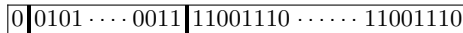
$$x \equiv \pm(1.a_1a_2 \dots a_m \dots)_2 \times 2^{(b_1b_2 \dots b_e)_2}$$



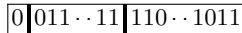
64-bit (Double)



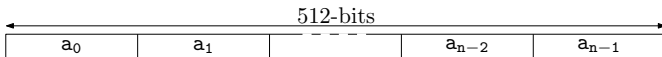
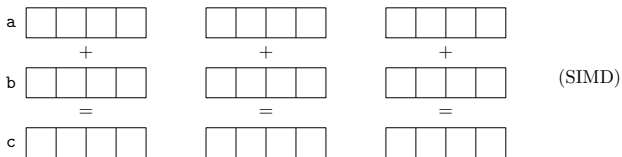
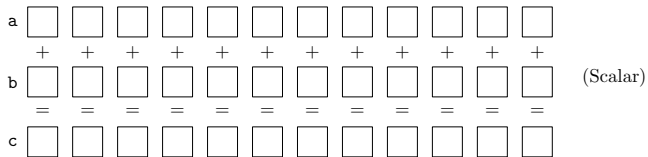
32-bit (Single)



16-bit (Half)

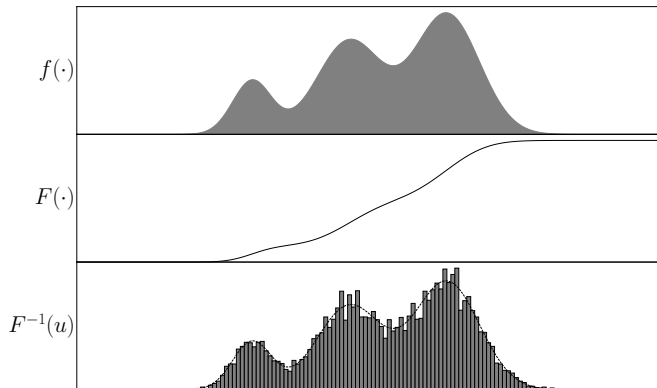


Vectorised (SIMD) operations

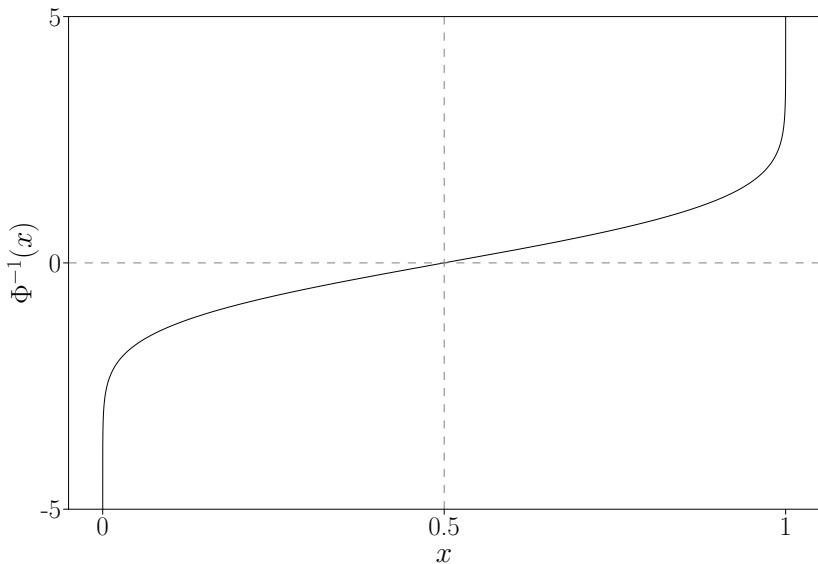


Mathematical situations of interest

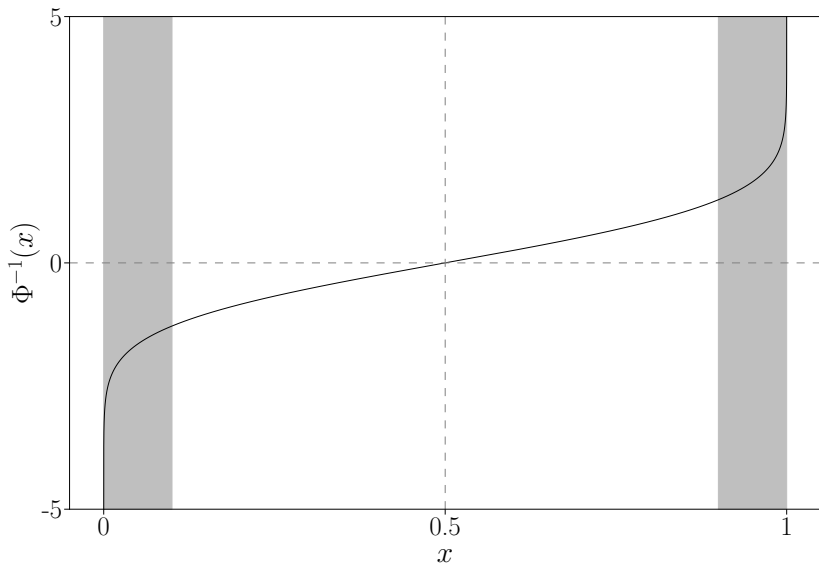
While there are many areas of mathematics which we could turn our attention to, we have been focusing on the numerical approximation of SDEs through Monte Carlo simulation. One of the primary areas where these concerns arise is in the simulation of random numbers from a distribution of interest, such as the Gaussian distribution.

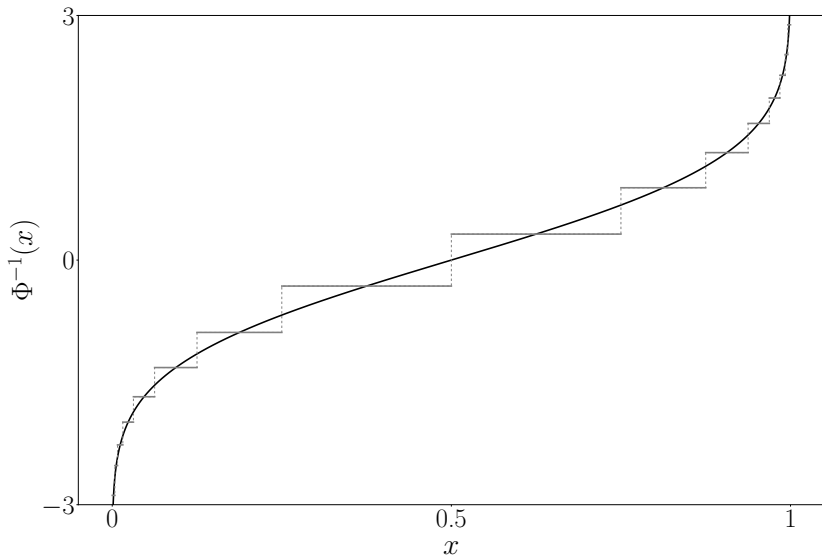


The inverse Gaussian CDF $\Phi^{-1}(\cdot)$



The inverse Gaussian CDF $\Phi^{-1}(\cdot)$





	Average speed (clock cycles)
Intel MKL	8
Lookup table	2

Accuracy and precision

HIGH PRECISION
HIGH ACCURACY



LOW PRECISION
HIGH ACCURACY



HIGH PRECISION
LOW ACCURACY



NO PRECISION
NO ACCURACY
NO SPECTATORS



J. HARRIS

$$\mathbb{E}(P) \approx \mathbb{E}(\hat{P}_L) = \mathbb{E}(\hat{P}_0) + \sum_{l=1}^L \mathbb{E}(\hat{P}_l - \hat{P}_{l-1}). \quad (1)$$

where $\mathbb{V}(\hat{P}_l - \hat{P}_{l-1}) \ll \mathbb{V}(\hat{P}_l)$ is achieved by **using the same underlying uniform random variables**.

Some possible multilevel constructions:

(Precision)

$$\mathbb{E}(\hat{P}_0) \rightarrow \mathbb{E}(\hat{P}_{32\text{-bit}}) + \mathbb{E}(\hat{P}_{64\text{-bit}} - \hat{P}_{32\text{-bit}})$$

$$\mathbb{E}(\hat{P}_0) \rightarrow \mathbb{E}(\hat{P}_{16\text{-bit}}) + \mathbb{E}(\hat{P}_{32\text{-bit}} - \hat{P}_{16\text{-bit}})$$

(Distribution)

$$\mathbb{E}(\hat{P}_0) \rightarrow \mathbb{E}(\hat{P}_{LT1024}) + \mathbb{E}(\hat{P}_Z - \hat{P}_{LT1024})$$

$$\mathbb{E}(\hat{P}_0) \rightarrow \mathbb{E}(\hat{P}_{\text{cubic}}) + \mathbb{E}(\hat{P}_Z - \hat{P}_{\text{cubic}})$$

We have the following SDE which we cannot solve:

$$dX_t = a(t, X_t) dt + b(t, X_t) dW_t, \quad (2)$$

where we wish to estimate some property which depends on the solution at the final time

$$\mathbb{E}(P(X_T))$$

(or possibly the entire duration).

$$\hat{X}_{n+1} = \hat{X}_n + a(\tau_n, \hat{X}_n)\delta + b(\tau_n, \hat{X}_n)\sqrt{\delta}Z_n \quad (3)$$

$$\hat{X}_{n+1} = \hat{X}_n + a(\tau_n, \hat{X}_n)\delta + b(\tau_n, \hat{X}_n)\sqrt{\delta}Z_n \quad (3)$$

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

$$\tilde{X}_{n+1} = \tilde{X}_n + a(\tau_n, \tilde{X}_n)\delta + b(\tau_n, \tilde{X}_n)\sqrt{\delta}\tilde{Z}_n \quad (4)$$

$$\check{X}_{n+1} = \check{X}_n + a(\tau_n, \check{X}_n)\delta + b(\tau_n, \check{X}_n)\sqrt{\delta}Z_n + \epsilon_n \quad (5)$$

$$\bar{X}_{n+1} = \bar{X}_n \oplus (\bar{a}(\bar{\tau}_n, \bar{X}_n) \otimes \bar{\delta} \oplus \bar{b}(\bar{\tau}_n, \bar{X}_n) \otimes (\sqrt{\bar{\delta}} \otimes \bar{Z}_n)) \quad (6)$$

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

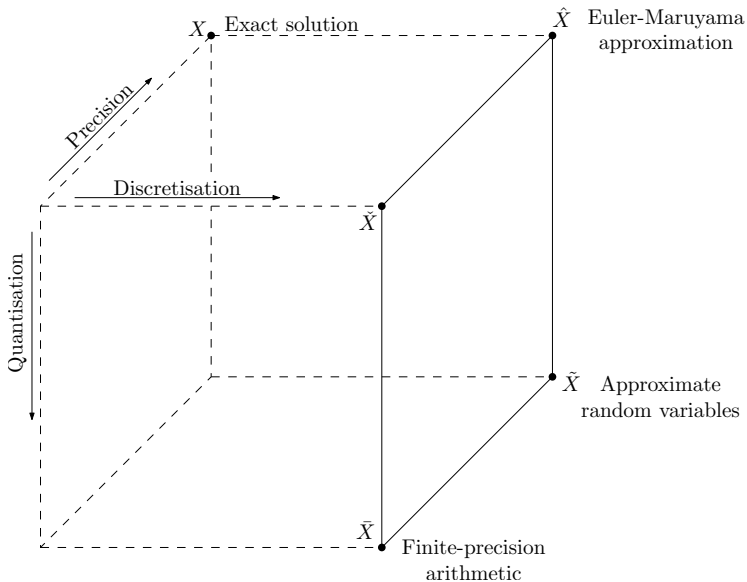
$$\tilde{X}_{n+1} = \tilde{X}_n + a(\tau_n, \tilde{X}_n)\delta + b(\tau_n, \tilde{X}_n)\sqrt{\delta}\tilde{Z}_n \quad (4)$$

$$\check{X}_{n+1} = \check{X}_n + a(\tau_n, \check{X}_n)\delta + b(\tau_n, \check{X}_n)\sqrt{\delta}Z_n + \epsilon_n \quad (5)$$

(Arciniega and Allen [1])

$$\bar{X}_{n+1} = \bar{X}_n \oplus (\bar{a}(\bar{\tau}_n, \bar{X}_n) \otimes \bar{\delta} \oplus \bar{b}(\bar{\tau}_n, \bar{X}_n) \otimes (\sqrt{\bar{\delta}} \otimes \bar{Z}_n)) \quad (6)$$

(Omland [2])



Theorem

Under Assumptions 11.1 to 11.6 where the constants do not depend on δ , and assuming $\mathbb{E}(\tilde{Z}) = 0$ and $\mathbb{V}(\tilde{Z}) < \infty$, then there exists a constant $K > 0$, independent of N , δ , and q , such that

$$\mathbb{E}\left(\sup_{n \leq N} |\hat{X}_n - \tilde{X}_n|^2\right) \leq K\mathbb{E}\left(|\tilde{Z} - Z|^2\right). \quad (7)$$

Theorem

Under Assumptions 11.1 to 11.6 where the constants do not depend on δ , we assume there exists a symmetric approximate normal distribution \tilde{Z} with a finite $(2p)$ -th moment such that $\mathbb{E}(\tilde{Z}^{2p-1}) = 0$ and $\mathbb{E}(|\tilde{Z}|^{2p}) < \infty$ for a finite integer $1 < p < \infty$. There exists a finite constant $0 < K_p < \infty$, independent of N , δ , and q , but dependent on p , such that

$$\mathbb{E}\left(\sup_{n \leq N} |\hat{X}_n - \tilde{X}_n|^{2p}\right) \leq K_p\mathbb{E}\left(|\tilde{Z} - Z|^{2p}\right). \quad (8)$$

Theorem

[1, Theorem 2.2] Under Assumptions 11.1 to 11.6 where the constants do not depend on δ , we use the adapted Euler-Maruyama scheme in (5). With i.i.d. errors $\epsilon_n \sim \mathcal{N}$ with $\mathbb{E}(\epsilon) = 0$ and $\mathbb{V}(\epsilon) \leq C\varrho^2$, where C is a some finite and strictly positive constant and $\varrho \propto 2^{-r}$, then there exists a constant $K > 0$, independent of N and r , such that

$$\mathbb{E} \left(\sup_{n \leq N} |\hat{X}_n - \check{X}_n|^2 \right)^{\frac{1}{2}} \leq K \sqrt{N} 2^{-r}. \quad (9)$$

$$\mathcal{R}(x \otimes y) = (x * y)(1 + \varepsilon), \quad (10)$$

$$|\bar{a}(t, x) - a(t, x)| \leq C2^{-r}(1 + |a(t, x)|), \quad (11)$$

Theorem

For all precisions $r \geq \lceil \log_2(N + 3) + 1 \rceil$, under Assumptions 11.1 to 11.4, assuming the roundoff model (10) and approximation model (11), and if the time domain $[t_0, T]$ is normalised to $[0, 1]$, and the times are partitioned on dyadic intervals such that $\delta = 2^{-k}$ for some $k \in \mathbb{Z}$, then the strong L^2 -error converges as

$$\mathbb{E} \left(\sup_{n \leq N} |\tilde{X}_n - \bar{X}_n|^2 \right)^{\frac{1}{2}} \leq CN2^{-r}, \quad (12)$$

where \bar{X} from (6) uses approximate random variables $\bar{Z} \sim \bar{\mu}_{\tilde{N}}$, and C is a finite strictly positive constant.

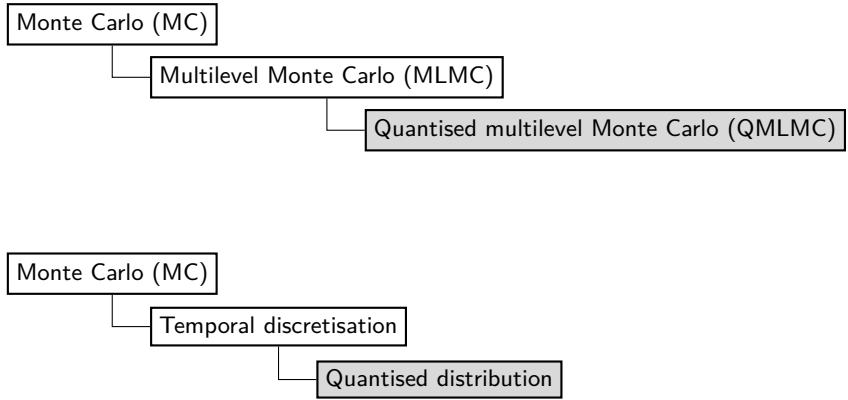
First we perform a multilevel decomposition with different **temporal discretisations**:

$$\mathbb{E}(\hat{P}_L) = \sum_l \mathbb{E}(\Delta_l \hat{P}_l). \quad (13)$$

Then for each of these terms we form a second multilevel decomposition using the **quantised distribution**:

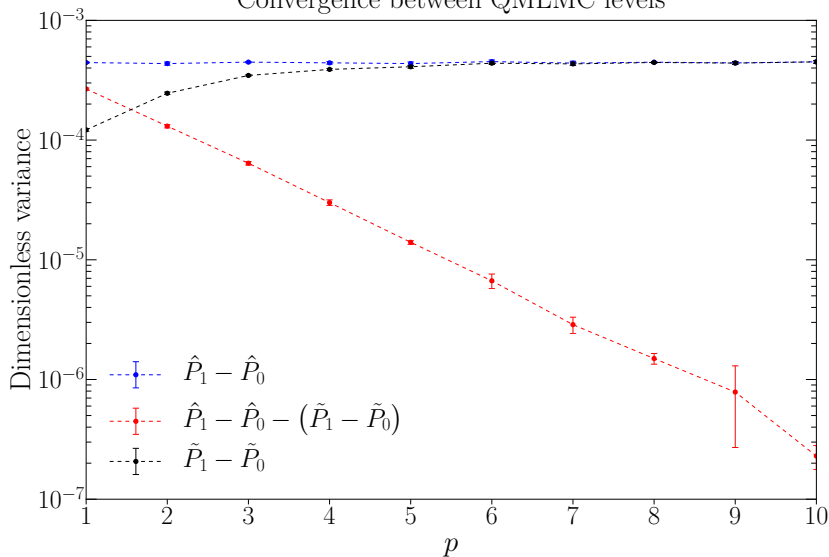
$$\mathbb{E}(\Delta_l \hat{P}_l) = \mathbb{E}(\Delta_l \tilde{P}_l) + \mathbb{E}(\Delta_l \hat{P}_l - \Delta_l \tilde{P}_l) \quad (14)$$

Nested multilevel Monte Carlo II



Nested multilevel Monte Carlo III

Convergence between QMLMC levels



For the exact Euler-Scheme there are several known results which we are trying to mirror/replicate. We have shown several strong convergence results from the quantised scheme to the discrete scheme, but would ultimately prefer to compare this directly with the exact solution.

$$\mathbb{E}(|\hat{X}_T^f - \hat{X}_T^c|^p) \preceq \delta^{\frac{p}{2}} \quad (15)$$

$$\mathbb{E}(|\tilde{X}_T^f - \tilde{X}_T^c|^p) \preceq \delta^{\frac{p}{2}} \quad (16)$$

$$\mathbb{E}(|(\hat{X}_T^f - \hat{X}_T^c) - (\tilde{X}_T^f - \tilde{X}_T^c)|^p) \preceq \delta^{\frac{p}{2}} \mathbb{E}(|Z - \tilde{Z}|^{p(1+\epsilon)})^{\frac{1}{1+\epsilon}} \quad (17)$$

$$\preceq \delta^{\frac{p}{2}} \mathbb{E}(|Z - \tilde{Z}|^p) \quad (18)$$

Running our MLMC estimator with a single quantisation level (1024 bins) for finely discretised paths gives the following average time per path:

Relative accuracy $\varepsilon = 10^{-3}$ Times per path 10^{-4} s	Memory intensive	Work intensive
Original MLMC	24.8	17.0
Quantised MLMC	13.2	3.99

Level	paths
Original	1 920 000
Quantised	1 980 000
Correction	14 000

Running our MLMC estimator with a single quantisation level (1024 bins) for finely discretised paths gives the following average time per path:

Relative accuracy $\varepsilon = 10^{-3}$ Times per path 10^{-4} s	Memory intensive	Work intensive
Original MLMC	24.8	17.0
Quantised MLMC	13.2	3.99

$\downarrow \times 2$

Level	paths
Original	1 920 000
Quantised	1 980 000
Correction	14 000

Running our MLMC estimator with a single quantisation level (1024 bins) for finely discretised paths gives the following average time per path:

Relative accuracy $\varepsilon = 10^{-3}$ Times per path 10^{-4} s	Memory intensive	Work intensive
Original MLMC	24.8	→17.0
Quantised MLMC	13.2	→3.99

Level	paths
Original	1 920 000
Quantised	1 980 000
Correction	14 000

Running our MLMC estimator with a single quantisation level (1024 bins) for finely discretised paths gives the following average time per path:

Relative accuracy $\varepsilon = 10^{-3}$ Times per path 10^{-4} s	Memory intensive	Work intensive
Original MLMC	24.8	17.0
Quantised MLMC	13.2	3.99

$\downarrow \times 4$

Level	paths
Original	1 920 000
Quantised	1 980 000
Correction	14 000

Errors from using a cheap proxy distribution can be **quantified** and controlled by the introduction of a nested multilevel Monte Carlo framework.

There is a degree of freedom in the construction of this proxy. Put the results in the low level **cache**, or use a very cheap (piece-wise) **polynomial**.

The resultant approximations from using such quantised distributions do **converge** under appropriate assumptions for several moments.

- [1] Armando Arciniega and Edward Allen. Rounding error in numerical solution of stochastic differential equations. *Stochastic Analysis and Applications*, 21(2):281–300, 2003.
- [2] Steffen Omland. Mixed precision multilevel Monte Carlo algorithms for reconfigurable computing systems, June 2016.

Assumption (Measurability)

The functions $a \equiv a(t, x)$ and $b \equiv b(t, x)$ are jointly (Lebesgue) \mathcal{L}^2 -measurable in $(t, x) \in [t_0, T] \times \mathbb{R}$.

Assumption (Lipschitz continuity)

There exists a constant $K > 0$ such that $|a(t, x) - a(t, y)| \leq K|x - y|$ for all $t \in [t_0, T]$ and $x, y \in \mathbb{R}$, and identically for b .

Assumption (Bounded growth)

There exists a constant $K > 0$ such that $|a(t, x)|^2 \leq K^2(1 + |x|^2)$ for all $t \in [t_0, T]$ and $x, y \in \mathbb{R}$, and identically for b .

Assumption (Measurable initial condition)

X_{t_0} is \mathcal{F}_{t_0} -measurable with $\mathbb{E}(X_{t_0}^2) < \infty$.

Furthermore, for our Euler-Maruyama schemes with approximate solution \hat{X}_T and uniform time increment δ we will further assume:

Assumption (Convergent initial value)

There exists a constant $K > 0$ such that $\mathbb{E}\left(|X_{t_0} - \hat{X}_{t_0}|^2\right) \leq K\delta$.

Assumption (Temporal Hölder continuity)

There exists a constant $K > 0$ such that $|a(t, x) - a(s, x)| \leq K(1+|x|)\sqrt{|t-s|}$ for all $s, t \in [t_0, T]$ and $x \in \mathbb{R}$, and identically for b .