



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

A Multilevel Monte Carlo Estimator for Matrix Multiplication

Citation for published version:

Wu, Y & Polydorides, N 2020, 'A Multilevel Monte Carlo Estimator for Matrix Multiplication', *SIAM Journal on Scientific Computing*, vol. 42, no. 5, pp. A2731-A2749. <https://doi.org/10.1137/19M125604X>

Digital Object Identifier (DOI):

[10.1137/19M125604X](https://doi.org/10.1137/19M125604X)

Link:

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

Document Version:

Peer reviewed version

Published In:

SIAM Journal on Scientific Computing

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 **A MULTILEVEL MONTE CARLO ESTIMATOR FOR MATRIX**
2 **MULTIPLICATION***

3 YUE WU [†] AND NICK POLYDORIDES [‡]

4 **Abstract.** Inspired by recent developments in multilevel Monte Carlo (MLMC) methods and
5 randomised sketching for linear algebra problems we propose a MLMC estimator for real-time pro-
6 cessing of matrix structured random data. Our algorithm is particularly effective in handling high-
7 dimensional inner products and matrix multiplication, and finds applications in computer vision and
8 large-scale supervised learning.

9 **Key words.** Sketching, inner product, matrix multiplication, real-time computing.

10 **AMS subject classifications.** 68W20.

11 **1. Introduction.** Randomised algorithms for matrix operations are in general
12 ‘pass-efficient’, and are primarily aimed at problems involving massive data sets that
13 are otherwise cumbersome to process with deterministic algorithms. Pass-efficient
14 implies that the algorithm necessitates only a very small number of passes through
15 the complete data set, but for the cases we consider here such a pass maybe turn
16 out to be impractical due to memory or time restrictions. In matrix multiplica-
17 tion for example, the BASICMATRIXMULTIPLICATION algorithm [3] is considered to
18 be the gold standard. Based on a probability assigned to the columns of a ma-
19 trix A , and respectively the rows of a matrix B , it approximates the product AB
20 through re-scaling the outer products of some sampled columns of A with the cor-
21 responding rows of B via a sampling-and-rescaling matrix operator. Variants of the
22 BASICMATRIXMULTIPLICATION algorithm were published in [4], [8], [13], exploiting
23 different types of information available on the elements of the matrices involved. In
24 particular, the algorithm in [4] addresses the case where the probability distribu-
25 tions of the elements are known a priori to devise an importance sampling strategy
26 based on BASICMATRIXMULTIPLICATION that minimizes the expected value of the
27 variance. The algorithm was shown to be effective when implemented with the opti-
28 mized sampling probabilities, particularly so in comparison to the estimators resulting
29 from uniform sampling. This result indeed extends BASICMATRIXMULTIPLICATION
30 to a random variable setting and can be applied to many query matching with in-
31 formation retrieval applications [4]. However, designing the optimized probabilities
32 relies exclusively on the knowledge of the probability distributions of the matrix el-
33 ements, which limits its applicability to the cases where such information is a priori
34 available. Conversely, it can be argued that BASICMATRIXMULTIPLICATION with
35 uniform probabilities becomes more appealing when dealing with real-time random
36 matrix multiplication tasks, where distributions change dynamically. In batch pro-
37 cessing for instance, the task at hand is to evaluate the expectation of the multipli-

*Submitted to the editors August 14, 2020.

Funding: NP and YW are grateful to EPSRC for funding this work through the project EP/R041431/1: ‘Randomness: a resource for real-time analytics’; YW is also funded by The Alan Turing Institute under the EPSRC grant EP/N510129/1 and by EPSRC through the project EP/S026347/1: ‘Unparameterised multi-modal data, high order signatures, and the mathematics of data science’.

[†]Mathematical Institute, University of Oxford, Oxford, UK; The Alan Turing Institute, London, UK; (yue.wu@maths.ox.ac.uk).

[‡]School of Engineering, University of Edinburgh, Edinburgh, UK; The Alan Turing Institute, London, UK; (n.polydorides@ed.ac.uk).

38 cation or indeed a functional of a matrix product at any given time, a formidable
 39 task in terms of the required speed and accuracy. To accelerate the time-dependent
 40 training of large-scale kernel machines for example, the evaluation of a kernel function
 41 is identified as and approximated through the expectation of a random inner product
 42 via some randomised feature map [10], [11]. In this case, coupling a standard Monte
 43 Carlo method (MC) and BASICMATRIXMULTIPLICATION with uniform probabilities
 44 may satisfy the speed specifications but compromise the accuracy of the result. A
 45 more prudent alternative is to employ a multilevel Monte Carlo method, similar to
 46 the one developed in [5] instead of MC.

47 MLMC was initially conceived for reducing the cost of computing the expected
 48 value of a financial derivative whose payoff depends upon the solution of a stochastic
 49 differential equation (SDE). The framework in [5] generalizes Kebaier’s approach in
 50 [9] to multiple levels, using a geometric sequence of different time step sizes. In do-
 51 ing so it reduces substantially the computational cost of MC by taking most of the
 52 samples on coarse grids at low cost and accuracy, and only a few samples on finer
 53 computationally expensive grids that lead to solutions of high accuracy. Over time,
 54 MLMC has grown in scope and found a wide range of applications in the broad area
 55 of SDEs, SPDEs, for stochastic reaction networks and inverse problems [12], while
 56 further variants have been developed in the form of multilevel quasi-Monte Carlo es-
 57 timators [6] and multilevel sequential Monte Carlo samplers [1]. For an overview on
 58 MLMC we refer the reader to the excellent survey [7]. Therein the author emphasizes
 59 that the multilevel theorem allows to use other estimators as long as they satisfy some
 60 specific conditions. This theorem lays the foundation for the algorithm proposed in
 61 this paper. A closely related work [2] considers the MLMC estimate for approximat-
 62 ing the mean field of a nonlinear PDE, providing a theoretical framework in separable
 63 Hilbert spaces. Although there is clearly no actual time stepsize in the matrix multi-
 64 plication context, we can draw an analogy between the term *time stepsize* in numerical
 65 analysis for differential equations and the term *the size of the sampled index set* in
 66 randomised linear algebra. As anticipated in a convergent MLMC scheme, the nu-
 67 merical estimation error shrinks with decreasing time stepsize. Similarly, due to the
 68 law of large numbers, increasing the size of index samples will decrease the expected
 69 squared Frobenius approximation error as shown in Lemma 4 of the seminal work [3].
 70 Therefore we claim that a random strategy for matrix multiplication with fewer index
 71 samples is analogous to using a “coarser grid” in the PDE setting. This observation
 72 is crucial to our construction of MLMC estimators for matrix multiplication.

73 In Section 2 below we begin by discussing the simpler case of calculating ‘on
 74 the fly’ the expectation of the inner product between large random vectors. We first
 75 consider the BASICMATRIXMULTIPLICATION algorithm with uniform probability and
 76 proceed to review the main results for the inner product from [4]. We then introduce
 77 the important quantities *base number* $M \in \mathbb{N}$ and *level size* $L \in \mathbb{N}$ based on which the
 78 MLMC estimator (c.f. (2.9) and (2.10)) is constructed via inner product approxima-
 79 tions with index sample sizes M^0, M^1, \dots, M^L . In this context, the approximation on
 80 the ‘finest grid’ corresponds to the inner product realization with M^L samples. Here
 81 we note the distinction between samples and indices, in that since we are sampling
 82 with replacement, taking $M^L > n$ samples does not imply sampling all n indices.
 83 Given that the variance of the approximated inner product is proportional to M^{-l}
 84 for $l \in \{1, \dots, L\}$ (c.f. Theorems 2.1 and Theorem 2.2), the complexity of the pro-
 85 posed MLMC estimator for a functional of the inner product conditioned on certain
 86 features of the underlying approximation can be treated similarly as the case $\beta = 1$ of
 87 Theorem 3.1 in [5]. This result is revisited in Theorem 2.2 where a comparison with

88 standard MC is attempted. Corollary 2.4 discusses the computational complexity of
 89 our MLMC estimator using Theorem 2.2. At the end of Section 2, we comment on
 90 the optimal choice of base number M following the reasoning in [5].

91 In Section 3 we extend our approach to matrix multiplication, adapting Theorems
 92 2.1 and 2.2 accordingly. It is worth mentioning that, because the approximation
 93 error (c.f. Theorems 3.1 and 3.2) is measured in expectation as a Frobenius norm,
 94 for the analysis the matrices are considered transformed in vector form prompting
 95 a new definition of ‘variance’ for the vectorized matrices denoted as \mathbb{V}_{\parallel} . Further,
 96 Theorem 3.3 discusses the complexity and Corollary 3.4 validates the complexity of
 97 the MLMC estimator for matrix multiplication. The implementation of our method
 98 is presented as Algorithm 3.2. Finally, in Section 4 we present two simple numerical
 99 experiments to illustrate the performance of the MLMC estimator in comparison with
 100 the standard MC one. By making appropriate choices for M and L parameters, the
 101 proposed MLMC estimator outperforms the MC estimator in terms of accuracy as
 102 well as speed and computational efficiency.

103 **2. Inner product.** We define \mathbf{T} as a countable collection of discrete time points
 104 and set $t \in \mathbf{T}$. Let $\mathbf{a}(t)$ and $\mathbf{b}(t)$ be two random vectors of length n , whose elements
 105 are drawn from some unknown, perhaps different, probability distributions, say $\mathbf{a}(t) \sim$
 106 $\mathcal{L}_{\mathbf{a}(t)}$ and $\mathbf{b}(t) \sim \mathcal{L}_{\mathbf{b}(t)}$. Here and throughout this paper, n is assumed to be extremely
 107 large such that evaluating the inner product of $\mathbf{a}(t)^T \mathbf{b}(t)$ is deemed impractical if at
 108 all possible. Consider that there is a need to compute $\mathbb{E}_{\mathbf{a}(t), \mathbf{b}(t)}[f(\mathbf{a}(t)^T \mathbf{b}(t))]$ on
 109 demand, at different times, where f is a Lipschitz function with Lipschitz constant C_f
 110 and $\mathbb{E}_{\mathbf{a}(t), \mathbf{b}(t)}$ is the expectation under $\mathcal{L}_{\mathbf{a}(t)}$ and $\mathcal{L}_{\mathbf{b}(t)}$. For the sake of notational
 111 simplicity, the argument (t) is suppressed in the notation but assumed implicitly in
 112 all of the quantities introduced above.

113 Indeed the task at hand consists of two main parts: approximating $\mathbf{a}^T \mathbf{b}$ in an
 114 efficient and accurate manner and approximating its expected value in the spirit of
 115 Monte Carlo methods. To tackle the first issue, the random sampling method for inner
 116 product estimation presents a viable option. Suppose there is a sampling distribution
 117 $\xi := \{\xi_j\}_{j=1}^n$ with $\sum_{j=1}^n \xi_j = 1$ such that each index $j \in [n]$, where $[n] := \{1, 2, \dots, n\}$,
 118 can be drawn with an assigned positive probability ξ_j . Further suppose we *fix* a ‘base’
 119 number $M \in \mathbb{N}$ and collect M^L , $L \in \mathbb{N}$, independent and identically distributed
 120 index samples as an index sequence (r_1, \dots, r_{M^L}) according to ξ . We shall refer to
 121 these collected M^L indices, or equivalently, the sequence (r_1, \dots, r_{M^L}) , as a sample
 122 *realisation*. Then denote by S_L the *sampling-and-rescaling matrix* of size $n \times M^L$ such
 123 that elements of \mathbf{a} and \mathbf{b} at the M^L index samples will be used for approximating the
 124 inner product of $\mathbf{a}^T \mathbf{b}$. That is,

$$125 \quad (2.1) \quad \widehat{\mathbf{a}^T \mathbf{b}} = \mathbf{a}^T S_L S_L^T \mathbf{b} = \frac{1}{M^L} \sum_{i=1}^{M^L} \frac{1}{\xi_{r_i}} \mathbf{a}_{r_i} \mathbf{b}_{r_i} := X_L(\xi),$$

127 where $X_L(\xi)$ denotes a scalar random variable that approximates the target $\mathbf{a}^T \mathbf{b}$
 128 using M^L samples from ξ , emphasizing its dependence on ξ . Previous research has
 129 shown that $X_L(\xi)$ is an unbiased estimator of $\mathbf{a}^T \mathbf{b}$ under the sampling distribution ξ .
 130 The performance of the approximation when the vector elements \mathbf{a} and \mathbf{b} are known
 131 only up to their distributions can be assessed through quantifying the variance of
 132 the estimator. The minimum variance is attained when sampling according to the
 133 distribution given by the following theorem from [4].

134 THEOREM 2.1. *If the vector elements \mathbf{a}_j and \mathbf{b}_j are independent random vari-*
 135 *ables, $j \in [n]$, with finite and nonzero moments $\mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{a}_j^2 \mathbf{b}_j^2]$, then the probability ξ^**
 136 *with elements*

$$137 \quad (2.2) \quad \xi_j^* = \frac{\sqrt{\mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{a}_j^2 \mathbf{b}_j^2]}}{\sum_{i=1}^n \sqrt{\mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{a}_i^2 \mathbf{b}_i^2]}},$$

139 *minimizes the expected value of the variance in (2.1), that is,*

$$140 \quad (2.3) \quad \min_{\xi} \mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{Var}[X_L(\xi)]] = \mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{Var}[X_L(\xi^*)]] := \frac{\mu}{M^L},$$

142 *where \mathbf{Var} is the variance under ξ and $\mu = \mathbb{E}_{\mathbf{a},\mathbf{b}}\left[\sum_{i=1}^n \frac{\mathbf{a}_i^2 \mathbf{b}_i^2}{\xi_i^*} - (\mathbf{a}^T \mathbf{b})^2\right]$.*

143 Sampling with ξ^* is clearly not practical when we have no knowledge about the
 144 distributions of \mathbf{a} and \mathbf{b} in advance, hence a plausible convenient alternative is to use
 145 a uniform probability over the index set

$$146 \quad (2.4) \quad \xi_j^u = \frac{1}{n}, \quad j \in [n],$$

148 with variance as follows.

149 THEOREM 2.2. [4] *Assume the same setting as in Theorem 2.1 but with probability*
 150 *ξ^u defined in (2.4), then the variance is*

$$151 \quad (2.5) \quad \mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{Var}[X_L(\xi^u)]] = \mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{Var}[X_L(\xi^*)]] + \frac{n\nu}{M^L} = \frac{n\nu + \mu}{M^L},$$

152 *where*

$$\nu = \sum_{i=1}^n \left(\sqrt{\mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{a}_i^2 \mathbf{b}_i^2]} - \frac{1}{n} \sum_{j=1}^n \sqrt{\mathbb{E}_{\mathbf{a},\mathbf{b}}[\mathbf{a}_j^2 \mathbf{b}_j^2]} \right)^2.$$

153 Typically one may consider approximating the expectation using a standard MC
 154 method that simulates $\mathbb{E}_{\mathbf{a},\mathbf{b}}[f(\mathbf{a}^T \mathbf{b})]$. In this instance, the quantity of interest, say P ,
 155 can then be estimated by (2.1) with a uniform probability (2.4) and MC as

$$156 \quad (2.6) \quad \mathbb{E}_{\mathbf{a},\mathbf{b}}[P] := \mathbb{E}_{\mathbf{a},\mathbf{b}}[f(\mathbf{a}^T \mathbf{b})] \approx \frac{1}{N} \sum_{k=1}^N f((\mathbf{a}^{(k)})^T \mathbf{b}^{(k)}) = \frac{1}{N} \sum_{k=1}^N f(X_L^{(k)}(\xi^u)) := \hat{P},$$

158 where N is the number of realisations for M^L many index samples or equivalently,
 159 an index sequence of length M^L . In this case the mean square error (MSE) for the
 160 estimate \hat{P} turns out to be

$$161 \quad (2.7) \quad \mathbb{E}[(\hat{P} - \mathbb{E}[P])^2] = \mathbb{E}[(\hat{P} - \mathbb{E}[\hat{P}])^2] + (\mathbb{E}[P] - \mathbb{E}[\hat{P}])^2$$

$$162 \quad = \mathbb{E}[(\hat{P} - \mathbb{E}[\hat{P}])^2] + (\mathbb{E}[f(\mathbf{a}^T \mathbf{b})] - \mathbb{E}[f(X_L(\xi^u))])^2,$$

163 where \mathbb{E} , and also \mathbb{V} that appears in the sequel, (without subscripts) denote respec-
 164 tively the expectation and the variance under $\mathcal{L}_{\mathbf{a}}$, $\mathcal{L}_{\mathbf{b}}$ and ξ^u . The last term in (2.7),
 165 for a fixed L , characterizes the bias and can be bounded by

$$166 \quad (\mathbb{E}[f(\mathbf{a}^T \mathbf{b})] - \mathbb{E}[f(X_L(\xi^u))])^2 \leq \mathbb{E}[(f(\mathbf{a}^T \mathbf{b}) - f(X_L(\xi^u)))]^2$$

$$167 \quad = \mathbb{E}[(f(\mathbf{a}^T \mathbf{b}) - f(\mathbf{a}^T S_L S_L^T \mathbf{b}))^2] \leq C_f^2 \mathbb{E}[|\mathbf{a}^T (\mathbf{I} - S_L S_L^T) \mathbf{b}|^2] \sim \mathcal{O}(M^{-L}),$$

169 where the first inequality comes from Jensen's inequality, that is $\mathbb{E}[X]^2 \leq \mathbb{E}[X^2]$ for
 170 arbitrary random variable X , the second inequality is due to the Lipschitz continuity
 171 of f and the last one due to (2.5). The first term in (2.7) is simply the variance from
 172 the MC simulation and can be bounded in terms of N as

$$\begin{aligned} \mathbb{E}[(\hat{P} - \mathbb{E}[\hat{P}])^2] &= \mathbb{V}[\hat{P}] = \frac{1}{N} \mathbb{V}[f(X_L(\xi^u))] \\ (2.8) \quad &\leq \frac{1}{N} \left(\mathbb{V}[f(X_L(\xi^u)) - f(\mathbf{a}^T \mathbf{b})]^{\frac{1}{2}} + \mathbb{V}[f(\mathbf{a}^T \mathbf{b})]^{\frac{1}{2}} \right)^2 \\ &\leq \frac{1}{N} \left(\frac{C_f}{M^{\frac{1}{2}}} (n\nu + \mu)^{\frac{1}{2}} + \mathbb{V}_{\mathbf{a}, \mathbf{b}}[f(\mathbf{a}^T \mathbf{b})]^{\frac{1}{2}} \right)^2 \sim \mathcal{O}(N^{-1}). \end{aligned}$$

175 Overall, as in [5], the MSE varies in terms of $\frac{1}{ML}$ and $\frac{1}{N}$. This is still true even if we
 176 sample based on the optimal sampling probability (2.2). Meanwhile, the complexity
 177 is in terms of NM^L , for an integer N to be determined.

178 Alternatively, it may be possible to obtain the same accuracy at a reduced compu-
 179 tational cost, by considering a multilevel MC simulation [5]. For $l \in [L] \cup \{0\}$ define
 180 as \hat{P}_l the approximation to $f(\mathbf{a}^T \mathbf{b})$ from M^l sampled indices. Further define \hat{Y}_l as
 181 an estimator of $\mathbb{E}[\hat{P}_l - \hat{P}_{l-1}]$ using N_l realizations with $l > 0$ and similarly \hat{Y}_0 the
 182 estimator of $\mathbb{E}[\hat{P}_0]$ using N_0 samples, that is

$$(2.9) \quad \hat{Y}_l := \frac{1}{N_l} \sum_{k=1}^{N_l} (\hat{P}_l^{(k)} - \hat{P}_{l-1}^{(k)}).$$

185 A key point to note is that both $\hat{P}_l^{(k)}$ and $\hat{P}_{l-1}^{(k)}$ emerge from the *same* realization, as
 186 we discuss in more detail when we describe our Algorithm 3.1. By the linear property
 187 of the expectation it follows immediately that

$$(2.10) \quad \mathbb{E}[\hat{P}_L] = \mathbb{E}[\hat{P}_0] + \sum_{l=1}^L \mathbb{E}[\hat{P}_l - \hat{P}_{l-1}] \approx \hat{Y}_0 + \sum_{l=1}^L \hat{Y}_l := \hat{Y},$$

190 where clearly $\mathbb{E}[\hat{P}_L] = \mathbb{E}[\hat{Y}]$. To investigate the performance of the proposed MLMC
 191 estimator \hat{Y} in (2.10) we compare the complexity of two estimators \hat{Y} and \hat{P} at the
 192 same accuracy level.

193 **THEOREM 2.3.** *Let \mathbf{a} and \mathbf{b} be two random vectors with length n drawn from*
 194 *different unknown distributions, that is $\mathbf{a} \sim \mathcal{L}_{\mathbf{a}}$ and $\mathbf{b} \sim \mathcal{L}_{\mathbf{b}}$, and let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a*
 195 *Lipschitz function with Lipschitz number C_f . Denote by P the term of interest as in*
 196 *(2.6), and define \hat{P}_l the corresponding approximation to $f(\mathbf{a}^T \mathbf{b})$ based on the sketched*
 197 *version of matrix multiplication via M^l many index samples like in (2.1).*

198 1. *If there exist independent estimators \hat{Y}_l as in (2.9) based on N_l Monte Carlo*
 199 *samples, and positive constants c_1, c_2, c_3 such that*

200 (a) $\mathbb{E}[\hat{P}_l - P] \leq c_1 M^{-\frac{1}{2}},$

201 (b) $\mathbb{V}[\hat{Y}_l] \leq c_2 N_l^{-1} M^{-l},$

202 (c) *the complexity of \hat{Y}_l , denoted by C_l , is bounded by $C_l \leq c_3 N_l M^l,$*

then there exists a positive constant c_4 such that for $\epsilon < e^{-1}$, there are values
 L and N_l for which the multilevel estimator $\hat{Y} = \sum_{l=0}^L \hat{Y}_l$ has an MSE $\mathbb{E}[(\hat{Y} -$
 $P)^2]$ with bound ϵ^2 , and computational complexity

$$C(\hat{Y}) := \sum_{l=0}^L C_l \leq c_4 \epsilon^{-2} (\log \epsilon)^2.$$

203 2. Furthermore, define the estimator based on the finest level L and N real-
 204 isations as in (2.6) with either the optimal sampling probability (2.2) (if
 205 tractable) or the uniform probability (2.4), and suppose

206 (a) the variance for \hat{P} is bounded by the same constant c_2 , i.e., $\mathbb{V}[\hat{P}] \leq$
 207 $c_2 N^{-1}$,

208 (b) the complexity for \hat{P} is bounded by the same constant c_3 , i.e., $C(\hat{P}) \leq$
 209 $c_3 N M^L$,

210 then at the same accuracy ϵ^2 , $C(\hat{P}) \leq c_6 \epsilon^{-4}$, which is much larger than the
 211 upper bound of $C(\hat{Y})$ when ϵ is sufficiently small.

212 *Proof.* 1. The proof is based on [5]. Accordingly, the MSE for \hat{Y} is

$$\begin{aligned} \mathbb{E}[(\mathbb{E}[P] - \hat{Y})^2] &= (\mathbb{E}[P] - \mathbb{E}[\hat{Y}])^2 + \mathbb{E}[(\hat{Y} - \mathbb{E}[\hat{Y}])^2] \\ &= (\mathbb{E}[P] - \mathbb{E}[\hat{P}_L])^2 + \mathbb{V}[\hat{Y}], \end{aligned}$$

216 where L is to be determined. If choosing the ceiling

$$(2.11) \quad L = \lceil \frac{\log(2c_1^2 \epsilon^{-2})}{\log M} \rceil,$$

219 then its bias component can be bounded via condition 1.(a)-(b) as

$$(\mathbb{E}[P] - \mathbb{E}[\hat{P}_L])^2 \leq c_1^2 M^{-L} \leq \frac{1}{2} \epsilon^2.$$

222 On the other hand, choosing

$$(2.12) \quad N_l = \lceil 2(L+1)c_2 \epsilon^{-2} M^{-l} \rceil$$

225 together with condition 1.(b) gives that

$$\begin{aligned} \mathbb{V}[\hat{Y}] &\leq \sum_{l=0}^L \mathbb{V}[\hat{Y}_l] \leq c_2 \sum_{l=0}^L N_l^{-1} M^{-l} \\ &\leq c_2 \sum_{l=0}^L (2(L+1)c_2 \epsilon^{-2} M^{-l})^{-1} M^{-l} \\ &= c_2 \sum_{l=0}^L \frac{\epsilon^2}{2(L+1)c_2} = \frac{1}{2} \epsilon^2. \end{aligned}$$

230 To bound the complexity C , let us first find the bound for L in terms of
 231 $\log \epsilon^{-1}$. Indeed, $L+1$, defined in (2.11) is bounded by

$$(2.13) \quad L+1 \leq \frac{2 \log(\epsilon^{-1})}{\log M} + \frac{\log(2c_1^2)}{\log M} + 2 \leq c_5 \log \epsilon^{-1},$$

234 where $c_5 = \frac{1 + (0 \vee \log(2c_1^2))}{\log M} + 2$ given that $\log \epsilon^{-1} > 1$ ($\epsilon \leq e^{-1}$). Besides, from
 235 (2.11) we can get an upper bound for M^{L-1} as

$$(2.14) \quad M^{L-1} \leq M^{\frac{\log(2c_1^2 \epsilon^{-2})}{\log M}} = e^{\log M \frac{\log(2c_1^2 \epsilon^{-2})}{\log M}} = 2c_1^2 \epsilon^{-2}.$$

238 Therefore the computational complexity C is bounded through

$$\begin{aligned}
239 \quad C &\leq c_3 \sum_{l=0}^L N_l M^l \leq c_3 \sum_{l=0}^L (2(L+1)c_2\epsilon^{-2}M^{-l} + 1)M^l \\
240 \quad &= c_3 \left(2(L+1)^2 c_2 \epsilon^{-2} + \frac{M^2 M^{L-1} - 1}{M-1} \right) \leq c_4 \epsilon^{-2} (\log \epsilon)^2, \\
241
\end{aligned}$$

242 where $c_4 = 2c_2 c_3 c_5^2 + \frac{2c_3 c_1^2 M^2}{M-1}$.

243 2. For both estimators \hat{Y} and \hat{P} , the bias is fixed for the same choice of L in
244 (2.11). Now let us choose an appropriate N such that $\mathbb{V}[\hat{P}] \leq \frac{1}{2}\epsilon^2$. Let
245 $N = \lceil 2c_2\epsilon^{-2} \rceil$ to meet the accuracy specification, and recall the upper bound
246 for M^{L-1} in (2.14). Then the complexity $C(\hat{P})$ is

$$247 \quad C(\hat{P}) \leq c_3 N M^L \leq c_3 (2c_2\epsilon^{-2} + 1) M^2 2c_1^2 \epsilon^{-2} \leq c_6 \epsilon^{-4},$$

248 where $c_6 = 2c_1^2 c_3 M^2 (2c_2 + e^{-2})$. □

250 The application of Theorem 2.3 relies on its conditions being verified. This is explored
251 in the form of the following corollary.

252 **COROLLARY 2.4.** *Assume the setting in Theorem 2.3 and choose a uniform sam-*
253 *pling distribution ξ^u as in (2.4). Then we have*

- 254 1. $c_1 = C_f^2(n\nu + \mu)$,
- 255 2. $c_2 = 2C_f^2(M+1)(n\nu + \mu) + 2\mathbb{V}_{\mathbf{a},\mathbf{b}}[P]$,
- 256 3. $c_3 = 1 + M^{-1}$.

257 *Proof.* 1. For any $l \in \mathbb{N} \cup \{0\}$ we have that

$$\begin{aligned}
258 \quad &(\mathbb{E}[f(\mathbf{a}^T \mathbf{b})] - \mathbb{E}[f(X_l(\xi^u))])^2 \leq \mathbb{E}[(f(\mathbf{a}^T \mathbf{b}) - f(X_l(\xi^u)))^2] \\
259 \quad &\leq C_f^2 \mathbb{E}[|\mathbf{a}^T (I - S_l S_l^T) \mathbf{b}|^2] \leq C_f^2 M^{-l} (n\nu + \mu),
\end{aligned}$$

261 where the last inequality holds because of (2.5).

262 2. For any $l > 0$ we have that

$$\begin{aligned}
263 \quad \mathbb{V}[\hat{P}_l - \hat{P}_{l-1}] &\leq (\mathbb{V}[\hat{P}_l - P]^{\frac{1}{2}} + \mathbb{V}[P - \hat{P}_{l-1}]^{\frac{1}{2}})^2 \\
264 \quad &\leq (\mathbb{E}[(\hat{P}_l - P)^2]^{\frac{1}{2}} + \mathbb{E}[(\hat{P}_{l-1} - P)^2]^{\frac{1}{2}})^2 \\
265 \quad &\leq C_f^2 (\mathbb{E}[|\mathbf{a}^T (I - S_l S_l^T) \mathbf{b}|^2]^{\frac{1}{2}} + \mathbb{E}[|\mathbf{a}^T (I - S_{l-1} S_{l-1}^T) \mathbf{b}|^2]^{\frac{1}{2}})^2 \\
266 \quad &\leq 2C_f^2 (M^{-l} + M^{-l+1}) (n\nu + \mu) \leq 2C_f^2 (M+1) (n\nu + \mu) M^{-l}.
\end{aligned}$$

268 For $l = 0$ we have that

$$\begin{aligned}
269 \quad \mathbb{V}[\hat{P}_0] &= \mathbb{V}[f(X_0(\xi))] \leq (\mathbb{V}[f(X_0(\xi)) - P]^{\frac{1}{2}} + \mathbb{V}[P]^{\frac{1}{2}})^2 \\
270 \quad &\leq (C_f \mathbb{E}[|\mathbf{a}^T (I - S_0 S_0^T) \mathbf{b}|^2]^{\frac{1}{2}} + \mathbb{V}_{\mathbf{a},\mathbf{b}}[P]^{\frac{1}{2}})^2 \\
271 \quad &\leq (C_f (n\nu + \mu)^{\frac{1}{2}} + \mathbb{V}_{\mathbf{a},\mathbf{b}}[P]^{\frac{1}{2}})^2 \\
272 \quad &\leq 2C_f^2 (n\nu + \mu) + 2\mathbb{V}_{\mathbf{a},\mathbf{b}}[P]. \\
273
\end{aligned}$$

274 Besides, from (2.8) we can see that $\mathbb{V}[\hat{P}]$ is bounded by the same c_2 .

3. For any $l > 0$ we can see easily the complexity is roughly

$$C_l \leq N^l(M^l + M^{l-1}) = (1 + M^{-1})N^l M^l,$$

while

$$C_0 \leq N^0 M^0 \leq (1 + M^{-1})N^0 M^0.$$

Besides, we have for the complexity of \hat{P} that

$$C(\hat{P}) \leq NM^L \leq (1 + M^{-1})NM^L.$$

Thus c_3 can be set as $1 + M^{-1}$. \square

Remark 2.5. Asymptotically as $l \rightarrow \infty$, we have that $\mathbb{E}[P - \hat{P}_l] \approx c_1 M^{-\frac{1}{2}}$, and hence

$$\mathbb{E}[\hat{P}_l - \hat{P}_{l-1}] \approx (\sqrt{M} - 1)c_1 M^{-\frac{1}{2}} \approx (\sqrt{M} - 1)\mathbb{E}[P - \hat{P}_l].$$

Similarly to the analysis in Section 4.2 of [5], this information can be used as an approximate bound: L can be set as the smallest integer such that

$$|\hat{Y}_L| < \frac{1}{\sqrt{2}}(\sqrt{M} - 1)\epsilon.$$

By doing this, we might achieve a bias bounded by $\frac{\epsilon^2}{2}$ without evaluating c_1 .

Remark 2.6 (Optimal N_l). To achieve a fixed variance, i.e., $\mathbb{V}[\hat{Y}] < \frac{1}{2}\epsilon$, the optimal N_l can be chosen as

$$N_l \approx \left\lceil 2\epsilon^{-2} \sqrt{V_l M^{-l}} \left(\sum_{j=0}^L \sqrt{V_l M^l} \right) \right\rceil,$$

where V_l denotes the variance of a single sample $\hat{P}_l - \hat{P}_{l-1}$. This result is simply an application of Section 1.3 of [7] or Eqn. (12) in [5] to the ‘stepsize’ M^{-l} . The estimation for N_l in (2.16) is conservative and may induce oversampling. In practice, some scaling factor might be introduced to avoid oversampling (see Section 4.1).

2.1. Optimal M . This part explores the methods in [5] in order to find an optimal M that reduces the computational complexity of the estimator even further. With c_2 given by Corollary 2.4, L and N_l given in the proof of Theorem 2.3, we can express the complexity of \hat{Y} in terms of M as

$$\begin{aligned} C(\hat{Y}) &\leq \sum_{l=0}^L C_l \approx \sum_{l=0}^L N_l (M^l + M^{l-1}) \stackrel{(2.12)}{\approx} \sum_{l=0}^L c_2 (L+1) (M^l + M^{l-1}) \epsilon^{-2} M^{-l} \\ &\stackrel{c_2}{\approx} \sum_{l=0}^L (L+1) (M+1)^2 M^{-1} \epsilon^{-2} = (L+1)^2 (M+1)^2 M^{-1} \epsilon^{-2} \\ &\stackrel{(2.11)}{\approx} M^{-1} (M+1)^2 \log(M)^{-2} \log(\epsilon)^2 \epsilon^{-2} = g(M) \log(\epsilon)^2 \epsilon^{-2}, \end{aligned}$$

where

$$g(M) := M^{-1} (M+1)^2 \log(M)^{-2}.$$

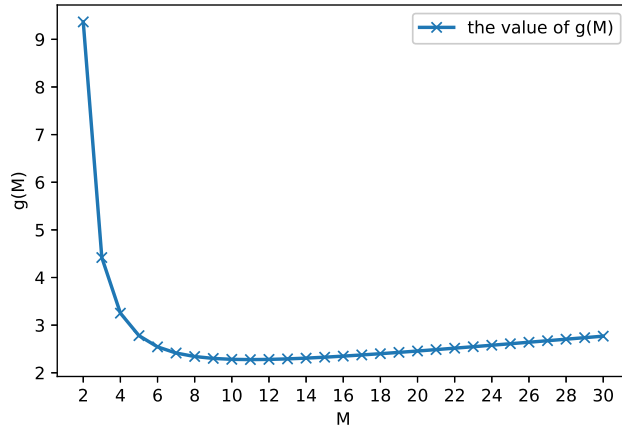


FIG. 1. The plot of the dominant complexity term $g(M)$ against the base number M , indicating the existence of an optimal M at the minimum point.

306 As illustrated in figure 1 where we plot $g(M)$ against M , $g(M)$ drops sharply for
 307 $M < 6$ and then starts growing slightly again after M going beyond 12. The minimum
 308 (optimum) is attained at $M = 11$, however from our experience using either $M = 10$
 309 or $M = 12$ does not make a significant difference. We remark that our definition of
 310 $g(M)$ in (2.17) differs somewhat from that used in [5], i.e. in the term $(M + 1)^2$, but
 311 this does not affect the general trend of $g(M)$ as described above. In the numerical
 312 experiments of Section 4.1, a choice of $M = 10$ is used as it was deemed appealing in
 313 terms of both the performance and time cost.

3. Matrix multiplication. We now extend our approach to matrix multiplication and thus we consider $A(t)$ and $B(t)$ to be two random matrices of size $m \times n$ and $n \times d$ respectively, $m, n, d \in \mathbb{N}$, drawn from different distributions, elementwise, in the sense $A(t) \sim \mathcal{L}_{A(t)}$ and $B(t) \sim \mathcal{L}_{B(t)}$, and again we suppress t in the notation as in Section 2 and assume that n is extremely large such that computing directly AB is prohibitively expensive. Recall that f is a Lipschitz function with Lipschitz constant C_f , and define $f^\circ(AB)$ the elementwise operator on AB , that is,

$$(f^\circ(AB))_{ik} = f((AB)_{ik}), \text{ for } i \in [m], \text{ and } k \in [d].$$

314 Once again consider that there is a need to compute $\mathbb{E}_{A,B}[f^\circ(AB)]$ where $\mathbb{E}_{A,B}$ is
 315 the expectation under \mathcal{L}_A and \mathcal{L}_B .

316 As in the inner product case, in order to simulate $\mathbb{E}_{A,B}[f^\circ(AB)]$ we first approxi-
 317 mate AB by random sampling (sketching) for matrix multiplication and then approxi-
 318 mate the expectation through a Monte Carlo method. Recall that $\xi := \{\xi_j\}_{j=1}^n$ with
 319 $\sum_{j=1}^n \xi_j = 1$ is a sampling probability such that an index $j \in [n]$ can be drawn with
 320 positive probability ξ_j and S_L a sampling-and-rescaling matrix of size $n \times M^L$ such
 321 that

$$(3.1) \quad \widehat{AB} = AS_L S_L^T B = \frac{1}{M^L} \sum_{i=1}^{M^L} \frac{1}{\xi_{r_i}} A_{:,r_i} B_{r_i,:} := Z_L(\xi),$$

324 where $A_{:,j}$ indicates the j th column of A and $B_{j,:}$ indicates the j th row of B , and
 325 $Z_L(\xi)$ denotes the matrix-valued random variable that approximates AB based on
 326 M^L indices sampled from ξ . It is easy to verify that $Z_L(\xi)$ is an unbiased estimator
 327 under the sampling distribution ξ . Besides, following arguments similar to those in
 328 the proof of Theorem 2.1 in [4] and Lemma 4 in [3], we can conclude that the minimum
 329 of the expected squared Frobenius error can be achieved by the following result.

330 **THEOREM 3.1.** *If the matrix elements A_{ij} and B_{jk} are independent random vari-*
 331 *ables, $i \in [m]$, $j \in [n]$ and $k \in [d]$, with finite and nonzero moments $\mathbb{E}_A[\|A_{:,j}\|_2^2]$ and*
 332 *$\mathbb{E}_B[\|B_{j,:}\|_2^2]$. Then the probability ξ^{**} , which is defined as*

$$333 \quad (3.2) \quad \xi_j^{**} = \frac{\sqrt{\mathbb{E}_A[\|A_{:,j}\|_2^2] \mathbb{E}_B[\|B_{j,:}\|_2^2]}}{\sum_{i=1}^n \sqrt{\mathbb{E}_A[\|A_{:,i}\|_2^2] \mathbb{E}_B[\|B_{i,:}\|_2^2]}}$$

334 *minimizes the expected value of the variance in (2.1), that is,*

$$335 \quad \min_{\xi} \mathbb{E}_{A,B} [\mathbf{E}[\|AB - Z_L(\xi)\|_F^2]] = \mathbb{E}_{A,B} [\mathbf{E}[\|AB - Z_L(\xi^{**})\|_F^2]]$$

$$336 \quad (3.3) \quad = \frac{1}{M^L} \left(\left(\sum_{j=1}^n \sqrt{\mathbb{E}_A[\|A_{:,j}\|_2^2] \mathbb{E}_B[\|B_{j,:}\|_2^2]} \right)^2 - \mathbb{E}_{A,B}[\|AB\|_F^2] \right) := \frac{\bar{\mu}}{M^L},$$

337 *where $\bar{\mu} = \left(\sum_{j=1}^n \sqrt{\mathbb{E}_A[\|A_{:,j}\|_2^2] \mathbb{E}_B[\|B_{j,:}\|_2^2]} \right)^2 - \mathbb{E}_{A,B}[\|AB\|_F^2]$, $\mathbf{E}[\cdot]$ denotes the ex-*
 338 *pectation under the distribution ξ , and $\mathbb{E}_{A,B}[\cdot]$ is the expectation with respect to the*
 339 *(element-wise) probabilities of A and B .*

341 The proof is omitted here as it is quite similar to the proof of Theorem 2.1 in [4].
 342 Besides, as discussed in Section 2, it is impractical to use ξ^* for random sampling. A
 343 simpler option would be to use a uniform probability ξ^u as defined in (2.4).

344 **THEOREM 3.2.** *Assume the same setting as in Theorem 3.1 but with probability*
 345 *ξ^u as defined in (2.4), then the expected squared Frobenius error is*

$$346 \quad (3.4) \quad \mathbb{E}_{A,B} [\mathbf{E}[\|Z_L(\xi^u) - AB\|_F^2]] = \mathbb{E}_{A,B} [\mathbf{E}[\|Z_L(\xi^{**}) - AB\|_F^2]] + \frac{n\bar{\nu}}{M^l} = \frac{n\bar{\nu} + \bar{\mu}}{M^l},$$

347 *where*

$$\bar{\nu} = \sum_{i=1}^n \left(\sqrt{\mathbb{E}_A[\|A_{:,i}\|_2^2] \mathbb{E}_B[\|B_{i,:}\|_2^2]} - \frac{1}{n} \sum_{j=1}^n \sqrt{\mathbb{E}_A[\|A_{:,j}\|_2^2] \mathbb{E}_B[\|B_{j,:}\|_2^2]} \right)^2.$$

348 The proof is omitted here as it is very similar to that of Theorem 2.3.

349 In this context, a quantity of interest P can be approximated with standard MC
 350 coupled to a random sampling method for matrix multiplication via either uniform
 351 probability (2.4) or the optimal probability (3.2) (if tractable)

(3.5)

$$352 \quad \mathbb{E}_{A,B}[P] := \mathbb{E}_{A,B}[f^\circ(AB)] \approx \frac{1}{N} \sum_{j=1}^N f(AS_L^{(j)}(S_L^{(j)})^T B) = \frac{1}{N} \sum_{j=1}^N f(Z_L^{(j)}(\xi)) := \hat{P},$$

353 *where N is the number of realisations for M^L many index samples. To consider the*
 354 *MSE for the estimate \hat{P} , we apply a matrix *vectorization*: for instance, if $A \in \mathbb{R}^{m \times n}$,*

$$355 \quad (3.6) \quad \text{vec}(A) = \text{vec}([A_{:,1} \quad \cdots \quad A_{:,n}]) = \begin{bmatrix} A_{:,1} \\ \vdots \\ A_{:,n} \end{bmatrix} \in \mathbb{R}^{mn},$$

358 is the column concatenation of A into a vector. Then the MSE would be

$$\begin{aligned}
359 \quad (3.7) \quad \mathbb{E}[\|\text{vec}(\hat{P} - \mathbb{E}[P])\|_2^2] &= \mathbb{E}[\|\text{vec}(\mathbb{E}[\hat{P}] - \mathbb{E}[P])\|_2^2] + \mathbb{E}[\|\text{vec}(\hat{P} - \mathbb{E}[\hat{P}])\|_2^2] \\
&= \|\text{vec}(\mathbb{E}[A(I - S_L S_L^T)B])\|_2^2 + \mathbb{E}[\|\text{vec}(\hat{P} - \mathbb{E}[\hat{P}])\|_2^2] \\
360 &= \|\text{vec}(\mathbb{E}[A(I - S_L S_L^T)B])\|_2^2 + \mathbb{V}_\parallel[\text{vec}(\hat{P})],
\end{aligned}$$

361 where \mathbb{E} is short for $\mathbb{E}_{A,B,\xi}$ and $\mathbb{V}_\parallel[\text{vec}(X)] := \mathbb{E}[\|\text{vec}(X - \mathbb{E}[X])\|_2^2]$ for any random
362 matrix X . Besides, it is easy to verify that

$$363 \quad (3.8) \quad \mathbb{V}_\parallel[X + Y]^{\frac{1}{2}} \leq \mathbb{V}_\parallel[X]^{\frac{1}{2}} + \mathbb{V}_\parallel[Y]^{\frac{1}{2}},$$

365 for any random vectors X and Y . Note that the variance of a vectorized random
366 matrix is indeed the variance of the random matrix in Frobenius norm. For example,

$$\begin{aligned}
367 \quad \mathbb{V}_\parallel[\text{vec}(\hat{P})] &= \mathbb{E}[\|\text{vec}(\hat{P} - \mathbb{E}[\hat{P}])\|_2^2] = \mathbb{E}\left[\sum_{h=1}^{md} \text{vec}(\hat{P} - \mathbb{E}[\hat{P}])_h^2\right] \\
368 &= \mathbb{E}\left[\sum_{i=1}^m \sum_{k=1}^d (\hat{P} - \mathbb{E}[\hat{P}])_{ik}^2\right] = \mathbb{E}[\|\hat{P} - \mathbb{E}[\hat{P}]\|_F^2]. \\
369
\end{aligned}$$

370 With these preliminaries let us now extend the approach of Section 3.1 to matrix
371 multiplication. For $l \in [L] \cup \{0\}$, define \hat{P}_l as the approximation to $f^\odot(AB)$ with M^l
372 many index samples. Recall that \hat{Y}_l is an estimator of $\mathbb{E}[\hat{P}_l - \hat{P}_{l-1}]$ using N_l realizations
373 with $l > 0$ and \hat{Y}_0 the respective estimator of $\mathbb{E}[\hat{P}_0]$ using N_0 samples, as defined in
374 (2.9). Eqn. (2.10) remains unchanged, from where we have that $\mathbb{E}[\hat{P}_L] = \mathbb{E}[\hat{Y}]$.

375 **THEOREM 3.3.** *Let A and B be two random matrices with sizes $m \times n$ and $n \times d$*
376 *respectively, drawn from different distributions, namely $A \sim \mathcal{L}_A$ and $B \sim \mathcal{L}_B$. Let*
377 *$f : \mathbb{R} \rightarrow \mathbb{R}$ be a Lipschitz function with Lipschitz number C_f . Denote by P the term*
378 *of interest as in (3.5). Define \hat{P}_l the corresponding approximation to $f^\odot(AB)$ based*
379 *on the sketched version of matrix multiplication via M^l many index samples like in*
380 *(3.1).*

381 1. *If there exist independent estimators \hat{Y}_l as in (2.9) based on N_l Monte Carlo*
382 *samples, and positive constants c_1, c_2, c_3 such that*

$$383 \quad \text{(a) } \|\text{vec}(\mathbb{E}[\hat{P}_l] - P)\|_2^2 \leq c_1^2 M^{-l},$$

$$384 \quad \text{(b) } \mathbb{V}_\parallel[\text{vec}(\hat{Y}_l)] \leq c_2 N_l^{-1} M^{-l},$$

385 *(c) the complexity of \hat{Y}_l , denoted by C_l , is bounded by $C_l \leq c_3 N_l M^l$,
then there exists a positive constant c_4 such that for $\epsilon < e^{-1}$, there are val-
ues L and N_l for which the multilevel estimator $\hat{Y} = \sum_{l=0}^L \hat{Y}_l$ has an MSE
 $\mathbb{E}[\|\text{vec}(\hat{Y} - \mathbb{E}[P])\|_2^2]$ with bound ϵ^2 , with computational complexity*

$$C(\hat{Y}) := \sum_{l=0}^L C_l \leq c_4 \epsilon^{-2} (\log \epsilon)^2.$$

386 2. *Furthermore, define the estimator based on the finest level L and N real-*
387 *izations as in (2.6) with either the uniform probability (2.4) or the optimal*
388 *probability (3.2) (if approachable). Suppose*

389 *(a) the variance for \hat{P} is bounded by the same constant c_2 , i.e., $\mathbb{V}_\parallel[\hat{P}] \leq$
390 $c_2 N^{-1}$,*

391 (b) the complexity for \hat{P} is bounded by the same constant c_3 , i.e., $C(\hat{P}) \leq$
 392 $c_3 N M^L$,
 393 then with the same accuracy ϵ^2 , $C(\hat{P}) \leq c_6 \epsilon^{-4}$ which is much larger than the
 394 bound of $C(\hat{Y})$.

395 The proof is similar to that of Theorem 2.3, expect from the decomposition of MSE,

$$396 \quad (3.9) \quad \mathbb{E}[\|\text{vec}(\hat{Y} - \mathbb{E}[P])\|_2^2] = \mathbb{E}[\|\text{vec}(\mathbb{E}[\hat{Y} - P])\|_2^2] + \mathbb{V}_\parallel[\text{vec}(\hat{Y})],$$

398 so we omit the proof. A more important issue is to verify our proposed MLMC
 399 estimator satisfies the conditions of Theorem 3.3.

400 COROLLARY 3.4. Assume the same setting in Theorem 3.3 via the sampling dis-
 401 tribution ξ^u in (2.4). Then we have

- 402 1. $c_1 = C_f^2(n\bar{\nu} + \mu)$,
- 403 2. $c_2 = 2C_f^2(M+1)(n\bar{\nu} + \mu) + 2\mathbb{V}_\parallel[f^\circ(AB)]$,
- 404 3. $c_3 = md(1 + M^{-1})$.

405 *Proof.* 1. For any $l \in \mathbb{N}$ we have that

$$406 \quad \begin{aligned} \|\text{vec}(\mathbb{E}[f^\circ(AB) - f^\circ(Z_l(\xi^u))])\|_2^2 &\leq \mathbb{E}[\|\text{vec}(f^\circ(AB) - f^\circ(Z_l(\xi^u)))\|_2^2] \\ 407 &= \mathbb{E}[\|f^\circ(AB) - f^\circ(Z_l(\xi^u))\|_F^2] \\ 408 &\leq C_f^2 \mathbb{E}[\|AB - Z_l(\xi^u)\|_F^2] \\ 409 &\leq C_f^2 M^{-l}(n\bar{\nu} + \mu), \end{aligned}$$

411 where the last inequality comes from Theorem 3.2.

412 2. For any $l > 0$ we have that

$$413 \quad \begin{aligned} &\mathbb{V}_\parallel[\text{vec}(\hat{P}_l - \hat{P}_{l-1})] \\ 414 &\leq \left(\mathbb{V}_\parallel[\text{vec}(\hat{P}_l - f^\circ(AB))]^{\frac{1}{2}} + \mathbb{V}_\parallel[\text{vec}(\hat{P}_{l-1} - f^\circ(AB))]^{\frac{1}{2}} \right)^2 \\ 415 &\leq \left(\mathbb{E}[\|\text{vec}(f^\circ(AB) - f^\circ(Z_l(\xi^u)))\|_2^2]^{\frac{1}{2}} \right. \\ 416 &\quad \left. + \mathbb{E}[\|\text{vec}(f^\circ(AB) - f^\circ(Z_{l-1}(\xi^u)))\|_2^2]^{\frac{1}{2}} \right)^2 \\ 417 &\leq 2C_f^2 (\mathbb{E}[\|AB - Z_l(\xi^u)\|_F^2]^{\frac{1}{2}} + \mathbb{E}[\|AB - Z_{l-1}(\xi^u)\|_F^2]^{\frac{1}{2}})^2 \\ 418 &\leq 2C_f^2 (M^{-l} + M^{-l+1})(n\bar{\nu} + \mu) \\ 419 &\leq 2C_f^2 (M+1)(n\bar{\nu} + \mu)M^{-l}. \end{aligned}$$

421 For $l = 0$ we have that

$$422 \quad \begin{aligned} &\mathbb{V}_\parallel[\text{vec}(\hat{P}_0)] = \mathbb{V}_\parallel[\text{vec}(f^\circ(Z_0(\xi^u)))] \\ 423 &\leq \left(\mathbb{V}_\parallel[\text{vec}(f^\circ(Z_0(\xi^u)) - f^\circ(AB))]^{\frac{1}{2}} + \mathbb{V}_\parallel[\text{vec}(f^\circ(AB))]^{\frac{1}{2}} \right)^2 \\ 424 &\leq 2C_f^2 \mathbb{E}[\|A(I - S_0 S_0^T)B\|_F^2] + 2\mathbb{V}_\parallel[\text{vec}(f^\circ(AB))] \\ 425 &\leq 2C_f^2(n\bar{\nu} + \mu) + 2\mathbb{V}_\parallel[\text{vec}(f^\circ(AB))]. \end{aligned}$$

427 Besides, it is easy to see that $\mathbb{V}_\parallel[\text{vec}(\hat{P})]$ can be bounded by the same c_2
 428 together with N^{-1} .

3. For any $l > 0$ we can see easily the complexity is roughly

$$C_l \leq mdN^l(M^l + M^{l-1}) = md(1 + M^{-1})N^l M^l,$$

while

$$C_0 \leq mdN^0 M^0 \leq md(1 + M^{-1})N^0 M^0.$$

Besides, we have for the complexity of \hat{P} that

$$C(\hat{P}) \leq mdNM^L \leq md(1 + M^{-1})NM^L.$$

Thus c_3 can be set as $md(1 + M^{-1})$. \square

Remark 3.5. Asymptotically as $l \rightarrow \infty$, we have that $\|\text{vec}(\mathbb{E}[P - \hat{P}_l])\|_2 \approx c_1 M^{-\frac{1}{2}}$, and hence

$$\|\text{vec}(\mathbb{E}[\hat{P}_l - \hat{P}_{l-1}])\|_2 \approx (\sqrt{M} + 1)c_1 M^{-\frac{1}{2}} \approx (\sqrt{M} + 1)\|\text{vec}(\mathbb{E}[P - \hat{P}_l])\|_2.$$

Similarly as in Section 4.2 of [5], this information can be used as an approximate bound: L can be set as the smallest integer such that

$$(3.10) \quad \|\text{vec}(\hat{Y}_L)\|_2 < \frac{1}{\sqrt{2}}(\sqrt{M} + 1)\epsilon.$$

By doing this, we might achieve a bias bounded by $\frac{\epsilon^2}{2}$ without evaluating c_1 .

Remark 3.6 (Optimal N_l). To achieve a fixed variance, i.e., $\mathbb{V}_{\parallel}[\text{vec}(\hat{Y})] < \frac{1}{2}\epsilon$, the optimal N_l can be chosen as

$$(3.11) \quad N_l \approx \left\lceil 2\epsilon^{-2} \sqrt{V_l M^{-l}} \left(\sum_{j=0}^L \sqrt{V_l M^l} \right) \right\rceil,$$

where V_l is the variance of the vectorized form of a single sample $\hat{P}_l - \hat{P}_{l-1}$ (recall the definition of the vectorized matrix variance right before (3.8)). This result is simply an application of Section 1.3 of [7] or Eqn. (12) in [5] to the ‘stepsize’ M^{-l} .

To choose the optimal value of M we argue as in Section 2.1, that is, $M = 11$ leads to the least computational complexity among between all choices of M . In the numerical experiments of Section 4.2, it turns out that $M = 10$ and $M = 6$ both yield acceptable approximations.

3.1. MLMC sketching algorithm. Based on the general discussion in the beginning of Section 2, Remark 3.5 and Remark 3.6, we propose an algorithm for estimating the matrix product based on MLMC method in Algorithm 3.2. The inner product case discussed in Section 2 can be treated as a special case. Algorithm 3.2 approximates $\mathbb{E}[f^\odot(AB)]$ through (2.10) under uniform probability (2.4), where the evaluation for each \hat{Y}_l in (2.10) is performed through function *level_estimation* described in Algorithm 3.1. To ensure the convergence, the choices of L and N_l with $l \in [L] \cup \{0\}$ are determined within Algorithm 3.2 using a *while* loop with one of the conditions given by Eqn. (3.10) in Remark 3.5¹. It is worth noticing that, while the

¹The condition for inner product is slightly different, see Eqn. (2.15) in Remark 2.5.

462 value L and therefore N_l for $l \in [L] \cup \{0\}$ are updated in the while loop (see Line 16
463 and Line 9), previous evaluations for \hat{Y}_l are reused in Line 13 for efficiency.

464 Although the outline in Algorithm 3.1 is simple to follow we draw the reader's
465 attention to Line 17 describing how $\hat{P}_l^{(k)}$ and $\hat{P}_{l-1}^{(k)}$ are computed through the com-
466 mon realization of M^l indices. Indeed, the procedure for getting $\hat{P}_l^{(k)}$ is by random
467 sampling as in (3.1) via the indices of a sample realization of size M^l under uni-
468 form probability, and likewise $\hat{P}_{l-1}^{(k)}$ via M^{l-1} of those M^l indices. That is, taking
469 (r_1, \dots, r_{M^l}) as a realization, then

$$470 \quad \hat{P}_l^{(k)} = f^\odot \left(\frac{n}{M^l} \sum_{j=1}^{M^l} A_{:,r_j}^{(k)} B_{r_j,:}^{(k)} \right), \quad \text{and} \quad \hat{P}_{l-1}^{(k)} = f^\odot \left(\frac{n}{M^{l-1}} \sum_{j=1}^{M^{l-1}} A_{:,r_{jM}}^{(k)} B_{r_{jM},:}^{(k)} \right).$$

472 Note that in practice the above computation can be further simplified by mapping
473 (r_1, \dots, r_{M^l}) into a set with non-repeated elements.

Algorithm 3.1 function *level_estimation*.

- 1: **Pre-defined:** \mathcal{L}_A and \mathcal{L}_B , the distributions of the targeted random matrices.
- 2: f , the targeted function; ξ^u , the uniform sampling distribution defined in (2.4).
- 3: **input:** l , the level size;
- 4: M , the base number;
- 5: N_l , the number of iterations.
- 6: **output:** \hat{Y}_l , the approximated version of $\mathbb{E}[\hat{P}_l - \hat{P}_{l-1}]$ for $l \neq 0$ or $\mathbb{E}[\hat{P}_l]$ for $l = 0$.
- 7: **initialization:** $\hat{Y}_l = 0$.
- 8: **if** $l = 0$ **then**
- 9: **for** $\ell = 1 \dots N_l$ **do**
- 10: get a pair of samples $A^{(\ell)}$ and $B^{(\ell)}$ from \mathcal{L}_A and \mathcal{L}_B ;
- 11: sample one index r from 1 to n according to ξ^u ;
- 12: set $\hat{Y}_l = \hat{Y}_l + \frac{1}{N_0} f^\odot(nA_{:,r}^{(\ell)} B_{r,:}^{(\ell)})$;
- 13: **else**
- 14: **for** $k = 1 \dots N_l$ **do**
- 15: get a pair of samples $A^{(k)}$ and $B^{(k)}$ from \mathcal{L}_A and \mathcal{L}_B ;
- 16: sample M^l many indices $(r_j)_{j=1}^{M^l}$ from 1 to n according to ξ^u ;
- 17: set

$$\hat{Y}_l = \hat{Y}_l + \frac{1}{N_l} \left(f^\odot \left(\frac{n}{M^l} \sum_{j=1}^{M^l} A_{:,r_j}^{(k)} B_{r_j,:}^{(k)} \right) - f^\odot \left(\frac{n}{M^{l-1}} \sum_{j=1}^{M^{l-1}} A_{:,r_{jM}}^{(k)} B_{r_{jM},:}^{(k)} \right) \right);$$

- 18: **return:** \hat{Y}_l .
-

474 **4. Numerical experiments.** In this part, we present some numerical experi-
475 ments designed to test the performance of the Algorithm 3.2 in comparison with a
476 standard MC method embedded with the optimal sampling distribution (see Theorem
477 2.1 and Theorem 3.1). Our experiments are implemented in Python (version 3.6.9)
478 with Numpy-based calculations being optimized under openBLAS [14] and executed
479 on a Linux cluster with two 14-core E5-2690 v4 Intel Xeon CPUs at 2.60GHz and
480 non-uniform memory allocation.

Algorithm 3.2 The MLMC estimator for $\mathbb{E}[f^\odot(AB)]$.

- 1: **Pre-defined:** \mathcal{L}_A and \mathcal{L}_B , the distributions of the targeted random vectors.
 - 2: f : the targeted function; ξ^u : the uniform distribution in (2.4).
 - 3: **input:** M , the base number;
 - 4: ϵ , the error tolerance.
 - 5: **output:** \hat{Y} , the approximated version of $\mathbb{E}[f^\odot(AB)]$.
 - 6: **initialization:** set $L = 0, t = 0$.
 - 7: **while** $L < 3$ or $\frac{\|\hat{Y}_{L-1}\|_F}{N_{L-1}^{(t-1)}} \geq \frac{1}{\sqrt{2}}(\sqrt{M} + 1)\epsilon$ **do**
 - 8: initialize $\hat{Y}_L = 0$;
 - 9: update V_l (defined in Remark 3.6) for all $l \in [L] \cup \{0\}$;
 - 10: calculate the optimal $N_l^{(t)}$ for all $l \in [L] \cup \{0\}$ through Eqn.(3.11);
 - 11: update $\hat{Y}_L = \hat{Y}_L + N_L^{(t)} \text{level_estimation}(L, M, N_L^{(t)})$;
 - 12: **if** $L > 0$ **then**
 - 13: **for** $l = 0 \dots L - 1$ **do**
 - 14: update $\hat{Y}_l = \hat{Y}_l + (N_l^{(t)} - N_l^{(t-1)}) \text{level_estimation}(l, M, N_l^{(t)} - N_l^{(t-1)})$;
 - 15: set $L = L + 1$ and $t = t + 1$;
 - 16: update $\hat{Y}_l = \hat{Y}_l / N_l^{(t-1)}$ for all $l \in [L - 1] \cup \{0\}$;
 - 17: **return:** $\sum_{l=0}^{L-1} \hat{Y}_l$.
-

481 **4.1. Example for the inner product.** Set $n = 10^4$ with $\mathbf{a}_j \sim \frac{j}{50}(0.4 - N(0, 1))$
482 and $\mathbf{b}_j \sim \cos(\text{Poi}(10) + 2\text{Exp}(1))\text{Bern}(0.05)$, $j \in [n]$, where $\text{Poi}(\lambda)$ is a Poisson random
483 variable with parameter λ , $\text{Exp}(\alpha)$ is an exponential random variable with parameter
484 α , and $\text{Bern}(\beta)$ is a Bernoulli random variable with success rate β . As the Bernoulli
485 random variable has low success rate we expect \mathbf{b} to be a sparse vector. In this
486 example we targeted function is set to $f(x) := \frac{1}{|x|H(x+0.4)+0.01}$, where $H(\cdot)$ is an
487 Heaviside step function. It is easy to see in this case f is highly nonlinear.

488 We test the Algorithm 3.2 for the inner product case with a parameter $M = 10$
489 and error tolerance $\epsilon = 0.1$, where the reference solution is obtained through direct
490 MC computation:

$$491 \quad (4.1) \quad \mathbb{E}[\mathbf{a}^T \mathbf{b}] \approx \frac{1}{\mathcal{N}_1} \sum_{j=1}^{\mathcal{N}_1} (\mathbf{a}^{(j)})^T \mathbf{b}^{(j)}, \quad \text{with } \mathcal{N}_1 = 10^5.$$

492 The value of L and the number of realizations at each level $l \leq L$, i.e., N_l with
493 $l \in [L] \cup \{0\}$ are tuned automatically by the algorithm itself. In our case $L = 3$ and
494 N_l is obtained through scaling (2.16) by $\frac{1}{20}$. This scaling factor $\frac{1}{20}$ is introduced to
495 prevent oversampling. Note that the scaling factor does not affect the trend of N_l .
496 Figure 2 illustrates the trend of the variance of each single path sample $\hat{P}_l - \hat{P}_{l-1}$
497 together with its corresponding N_l . From there it is easy to see that there is a
498 clear decay in variance with respect to l from $l = 1$, which results in the nearly
499 polynomial decay in the number of realizations N_l . For comparison we also perform
500 a standard MC simulation of the same M and L under *optimal sampling* (Theorem
501 2.1) with a number of repetitions chosen to maintain roughly the same accuracy level
502 (convergence). The results obtained are tabulated in Table 1.

503 From Table 1, the MLMC estimator using ξ^u in general outperforms the MC one
504 using ξ^* in terms of the elapsed time. Though MC using ξ^u provides an approxima-
505 tion that doubles the accuracy of MLMC, its computational time is about six times

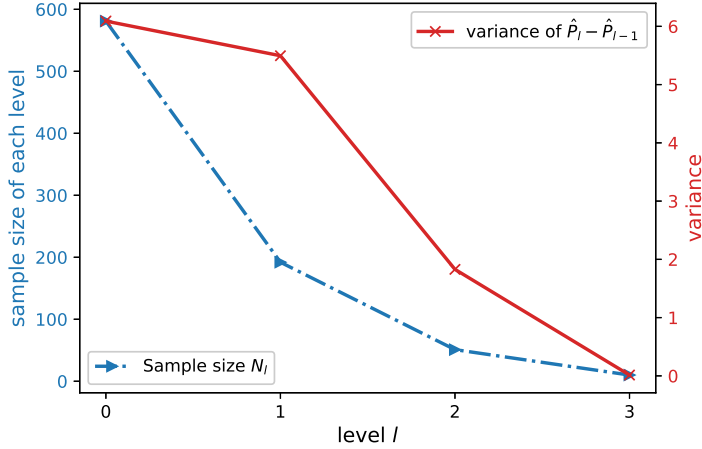


FIG. 2. Inner product case: plot of the variance of a single realization $\hat{P}_l - \hat{P}_{l-1}$ up to $l = L$ (red solid line), and its corresponding number of realizations for each l up to $l = L$ (blue dashed line): for $l = 0$, the variance of a single realization $\hat{P}_l - \hat{P}_{l-1}$ is indeed the variance of a single realization \hat{P}_0 .

| (M, L) | MLMC using ξ^u | | | MC using ξ^* | | | directMC |
|----------|--------------------|-------|-----------|------------------|-------|-----------|-----------|
| | AE | RE | time cost | AE | RE | time cost | time cost |
| (10, 3) | 0.002 | 0.079 | 0.047 s | 0.001 | 0.041 | 0.274 s | 0.423 s |

TABLE 1

Numerical results from the implementation of our method on approximating the inner product. These include records of the relative errors (RE), absolute errors (AE) and computational times for Algorithm 3.2 under $M = 10$ and its corresponding L . For comparison we provide also the results from standard MC (2.6) with optimal sampling distribution ξ^* (2.2) based on the finest level L , and time cost for getting reference solution through direct MC (4.1).

506 longer. The computation times of both estimators are less than that of the directMC
 507 for getting the reference solution, which illustrates the advantage of our proposed
 508 estimator in practice.

509 **4.2. Example for the matrix multiplication.** In the matrix multiplication
 510 case we consider a setup with $n = 10^4$, $m = d = 10^3$ using $A_{ij} \sim g_1\left(\frac{j}{10^4}(0.5 -\right.$
 511 $\left. N(0, 1))\right)$, where $g_1(x) := \sin(x) + N(0, 1)x$, and $B_{jk} \sim g_2(\text{Poi}(2))\text{Bern}(0.2)$, where
 512 $g_2(x) := \cos(x)H(5 - x)$ for $i \in [m], j \in [n]$ and $k \in [d]$. Like before, $\text{Poi}(\lambda)$ denotes a
 513 Poisson random variable with parameter λ and $\text{Bern}(\beta)$ is a Bernoulli random variable
 514 with success rate β . The targeted function is chosen to be $f(x) := |x|H(2 - x)$, where
 515 $H(\cdot)$ is an Heaviside step function.

516 Similar to the inner product example, we run Algorithm 3.2 for the matrix product
 517 with base number $M = 10$ and the error tolerance $\epsilon = 0.1$. The reference solution is
 518 computed as

$$519 \quad (4.2) \quad \mathbb{E}[AB] \approx \frac{1}{\mathcal{N}_2} \sum_{j=1}^{\mathcal{N}_2} A^{(j)} B^{(j)}, \quad \text{with } \mathcal{N}_2 = 10^5.$$

520 Algorithm 3.2 automatically chooses $L = 5$. Though M^L is now larger than n , this

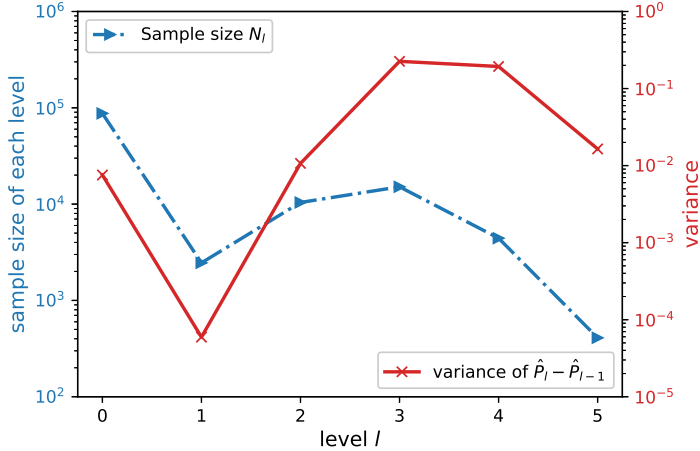


FIG. 3. Matrix multiplication case: Plot of the variance of a single realization $\hat{P}_l - \hat{P}_{l-1}$ up to $l = L$ (red solid line), and its corresponding number of realizations for each l up to $l = L$ (blue dashed line): for $l = 0$, the variance of a single realization $\hat{P}_l - \hat{P}_{l-1}$ is indeed the variance of a single realization \hat{P}_0 .

521 does not imply sampling all the columns of A . Besides, the number of realizations
 522 generated at the finest level L is very small, which does not affect the total perfor-
 523 mance. Meanwhile the variance of each single path sample $\hat{P}_l - \hat{P}_{l-1}$ together with its
 524 corresponding N_l is directly obtained through (3.11). Figure 3 illustrates the trend
 525 of the variance of $\hat{P}_l - \hat{P}_{l-1}$ and N_l up to $l = L$. It can be noted from Figure 3
 526 that from $l = 3$ the variance curve begins to decay, while the apparent low values
 527 in variance from $l = 0$ to $l = 2$ are due to the nature of randomised sketching of
 528 matrix multiplication. For example, to get an approximation $Z_l(\xi^u)$ of AB through
 529 (3.1) with $l = 0$, only one column of A and the corresponding row of B are selected
 530 and multiplied. This is guaranteed to decrease the variance of $Z_0(\xi^u)$. We can also
 531 observe this phenomenon from the inner product case as depicted in Figure 2, where
 532 the variance of $l = 0$ is only slightly bigger than the one of $l = 1$. The curve in Figure
 533 3 also indicates that our proposed estimator will be more efficient in super-large-scale
 534 matrix application, where the variance decay speeds up for higher level $l \gg 5$.

535 To compare performance, a standard MC simulation of the same M and L , for-
 536 mulated in (3.5), is implemented with optimal sampling distribution ξ^{**} (Theorem
 537 3.1) with the number of repetitions chosen to maintain roughly the same accuracy
 538 level. As matrix B is a very sparse matrix, optimal probability defined in Eqn. 3.3
 539 might have sparse or very small entries. Therefore even M^L is now larger than n ,
 540 the probability that all the columns of A are sampled to obtain an approximation is
 541 pretty small. The results obtained are recorded in Table 2. The MLMC estimator
 542 using ξ^u in general outperforms the MC one using ξ^{**} in terms of the elapsed time.
 543 Meanwhile, the computational times for MC using ξ^{**} are admittedly very large, tak-
 544 ing three times longer than the directMC. This is mainly due to the choice of the
 545 high level $L = 5$ compared to the matrix size. On the other hand, it is reasonable
 546 to anticipate that the MLMC method under the approximated optimal probability
 547 instead of the uniform one, would lead to a drastic improvement of the efficiency of
 548 the approximation beyond what has been demonstrated in this work.

| (M, L) | MLMC using ξ^u | | | MC using ξ^{**} | | | directMC |
|----------|--------------------|-------|-----------|---------------------|-------|-----------|-----------|
| | AE | RE | time cost | AE | RE | time cost | time cost |
| (10, 5) | 0.088 | 0.006 | 2.240 s | 0.069 | 0.005 | 75.173 s | 25.561 s |

TABLE 2

Numerical results from the implementation of our method on approximating the matrix product. These include records of the absolute errors in Frobenius norm (AE), the relative errors (RE) and computational times for Algorithm 3.2 under $M = 10$ and its corresponding L . For the sake of comparison we provide also the results from a standard MC simulation (3.5) based on the finest level L and sampling distribution ξ^{**} (3.2), and the time cost for getting the reference solution through direct MC (4.2)

549 **5. Conclusions.** We presented a new approach for computing arbitrary vector
550 and matrix products ‘on-the-fly’ that combines ideas from sketching in randomized
551 linear algebra and multilevel Monte Carlo approaches for estimating high-dimensional
552 integrals. Our approach is simple to implement and, subject to optimizing some
553 algorithmic parameters, it outperforms the standard Monte Carlo in both in terms of
554 the accuracy and the time required for computing the estimator.

555 **Acknowledgements.** The authors are grateful to EPSRC for funding this work
556 through the project EP/R041431/1, titled ‘Randomness: a resource for real-time an-
557 alytics’; YW is also funded by The Alan Turing Institute under the EPSRC grant
558 EP/N510129/1 and by EPSRC through the project EP/S026347/1, titled ‘Unpa-
559 rameterised multi-modal data, high order signatures, and the mathematics of data
560 science’.

561

REFERENCES

- 562 [1] Beskos, A., Jasra, A., Law, K., Tempone, R., and Zhou, Y. (2017). *Multilevel sequential Monte*
563 *Carlo samplers*. Stochastic Processes and their Applications, 127(5), 1417-1440.
564 [2] Bierig, C. and Chernov, A. (2015). *Convergence analysis of multilevel Monte Carlo variance*
565 *estimators and application for random obstacle problems*. Numerische Mathematik, 130(4),
566 579-613.
567 [3] Drineas, P., Kannan R, and Mahoney, W. M. (2006). *Fast Monte Carlo algorithms for matrices*
568 *I: Approximating matrix multiplication*, SIAM J. Comput. 36(1), 132-157.
569 [4] Eriksson-Bique, S., Solbrig, M., Stefanelli, M., Warkentin, S., Abbey, R., Ipsen, I.C.F. (2011).
570 *Importance sampling for a Monte Carlo matrix multiplication algorithm, with application*
571 *to information retrieval*, SIAM J. Comput., 33, 16891706.
572 [5] Giles, M. B. (2008). *Multilevel monte carlo path simulation*. Operations Research 56.3, 607-617.
573 [6] Giles, M. B. and Waterhouse, B. J. (2009). *Multilevel quasi-Monte Carlo path simulation*.
574 *Advanced Financial Modelling, Radon Series on Computational and Applied Mathematics*
575 8, 165-181.
576 [7] Giles, M. B. (2015) *Multilevel monte carlo methods*. Acta Numerica 24, 259-328.
577 [8] Holodnak, J., and Ipsen, I. (2015). *Randomized approximation of the gram matrix: Exact*
578 *computation and probabilistic bounds*. SIAM Journal on Matrix Analysis and Applications
579 36.1: 110-137.
580 [9] Kebaier, A. (2005). *Statistical Romberg extrapolation: a new variance reduction method and*
581 *applications to option pricing*. The Annals of Applied Probability 15.4 (2005): 2681-2705.
582 [10] Kar, P., and Karnick, H. (2012). *Random feature maps for dot product kernels*. In Artificial
583 *Intelligence and Statistics*, 583-591.
584 [11] Rahimi, A., and Recht, B. (2008). *Random features for large-scale kernel machines*. In Advances
585 *in neural information processing systems*, 1177-1184.
586 [12] Teckentrup, A.L., Scheichl, R., Giles, M. B., and Ullmann E. (2013). *Further analysis of multi-*
587 *level Monte Carlo methods for elliptic PDEs with random coefficients*. Numerische Math-
588 *ematik*, 125(3), 569-600.

- 589 [13] Wu, Y. (2018). *A Note on Random Sampling for Matrix Multiplication*. arXiv preprint,
590 arXiv:1811.11237.
- 591 [14] Zhang, X., Wang, Q. and Chothia, Z. (2012). *openBLAS*. <http://xianyi.github.io/OpenBLAS>,
592 88.