# Deflated Multigrid Multilevel Monte Carlo

**Andreas Frommer**[a] **and Gustavo Ramirez-Hidalgo**[a,*]

[a]*Department of Mathematics, Bergische Universität Wuppertal, 42097 Wuppertal, Germany*

*E-mail:* g.ramirez@math.uni-wuppertal.de

In lattice QCD, the trace of the inverse of the discretized Dirac operator appears in the disconnected fermion loop contribution to an observable. As simulation methods get more and more precise, these contributions become increasingly important. Hence, we consider here the problem of computing the trace $\mathrm{tr}(D^{-1})$, with $D$ the Dirac operator. The Hutchinson method, which is very frequently used to stochastically estimate the trace of a function of a matrix, approximates the trace as the average over estimates of the form $x^H D^{-1} x$, with the entries of the vector $x$ following a certain probability distribution. For $N$ samples, the accuracy is $O(1/\sqrt{N})$. In recent work, we have introduced multigrid multilevel Monte Carlo: having a multigrid hierarchy with operators $D_\ell$, $P_\ell$ and $R_\ell$, for level $\ell$, we can rewrite the trace $\mathrm{tr}(D^{-1})$ via a telescopic sum with difference-levels, written in terms of the aforementioned operators and with a reduced variance. We have seen significant reductions in the variance and the total work with respect to exactly deflated Hutchinson. In this work, we explore the use of exact deflation in combination with the multigrid multilevel Monte Carlo method, and demonstrate how this leads to both algorithmic and computational gains.

---

*Speaker

## 1. Introduction

In lattice quantum chromodynamics (LQCD), the extraction of some observables requires the computation of matrix traces [1] of the form $\text{tr}(f(D))$ with $D$ the Dirac operator on the lattice, and $f(D) = (\Gamma D)^{-1}$ with $\Gamma$ some Dirac structure e.g. $\Gamma = \Gamma_5$. This appears, for example, in the calculation of disconnected diagrams [1].

There are deterministic algorithms for the approximate computation of such traces, among which falls hierarchical probing [2, 3]. From the stochastic side, on which we focus in this work, a well known method is the use of the Hutchinson estimator [4]

$$\text{tr}(f(A)) \approx \frac{1}{N} \sum_{i=1}^{N} (x^{(i)})^H f(D) x^{(i)}. \tag{1}$$

The entries $x_j$ in the vectors $x^{(i)}$ in eq. (1) are random identically and independently distributed (i.i.d.) with expected values $\mathbb{E}[x_i] = 0$ and $\mathbb{E}[x_i x_j] = \delta_{ij}$. The variance $\mathbb{V}[x^H f(D) x]$, in terms of the entries of the matrix $f(D)$, is determined by the probability distribution (p.d.f.) used for drawing the components of the vector $x$.

With $\mathbb{V}(x^H f(D) x)$ having a theoretical value that depends on $f(D)$ and the p.d.f. used in drawing $x$, the variance of the Hutchinson estimator in eq. (1) decreases as $\frac{1}{\sqrt{N}}$. When the accuracy required in the computation of $\text{tr}(f(D))$ is quite high, the Hutchinson method becomes too costly. There are multiple variance reduction techniques that can be applied to the Hutchinson estimator, in particular deflation [5], which we discuss in section 2.

A new variance reduction technique recently developed in [6], which we briefly describe in section 3, uses a multilevel approach via a multigrid hierarchy in order to "distribute" the variance over different multigrid levels, and with this offloading some (and in some cases most) of the computational work to coarser levels.

In this work, we discuss the use use of exact deflation on the multigrid multilevel Monte Carlo method from [6]. We focus on $f(D) = D^{-1}$, with $D$ the Wilson-Schwinger operator coming from a discretization of the (1+1)-dimensional Schwinger model on the lattice.

## 2. Deflated Hutchinson

If $\mathbb{Z}_2$ noise is used for drawing the random vectors used in the Hutchinson estimator, then [7, 8]

$$\mathbb{V}(x^H D^{-1} x) = \frac{1}{2} \|\text{offdiag}(D^{-1} + D^{-T})\|_F^2. \tag{2}$$

Now, for a given matrix $A \in \mathbb{C}^{n \times n}$, there is a connection between the Frobenius norm of $A$ and its singular values

$$\|A\|_F^2 = \sum_{i=1}^{n} \sigma_i^2 \implies \|\text{offdiag}(A)\|_F^2 = \sum_{i=1}^{n} \sigma_i^2 - \sum_{i=1}^{n} |A_{ii}|^2. \tag{3}$$

In LQCD simulations, where the smallest singular values of $D$ are typically very small—and those of $D^{-1}$ are thus large—, the last term in eq. (3) can be discarded, and we see that the smallest

singular modes of $D$ are the ones dominating the variance. Deflating those modes can then lead to a significant reduction in the variance of the Hutchinson estimator.

Exactly deflated Hutchinson [5] has been extensively used in LQCD, e.g. [9–11]. With the singular value decomposition (SVD) [12] of the Dirac operator, $D = U\Sigma V^H$ and its inverse $D^{-1} = V\Sigma^{-1}U^H$, exact deflation is done via the orthogonal projector

$$\Pi = U_k U_k^H, \tag{4}$$

where the columns of $U_k$ contain the $k$ largest right singular vectors of $D^{-1}$, i.e., the $k$ smallest left singular vectors of $D$. The trace $\mathrm{tr}(D^{-1})$ can then be split as

$$\mathrm{tr}(D^{-1}) = \mathrm{tr}(D^{-1}(I - \Pi)) + \mathrm{tr}(D^{-1}\Pi). \tag{5}$$

From the discussion following eq. (2), the first term in eq. (5) will have a reduced variance, so we can compute that term via the Hutchinson estimator investing less estimates. For the second term, we can write

$$\mathrm{tr}(D^{-1}\Pi) = \mathrm{tr}(D^{-1}U_k U_k^H) = \mathrm{tr}(U_k^H D^{-1}U_k) = \mathrm{tr}(U_k^H V_k \Sigma_k^{-1}), \tag{6}$$

where the last step requires a pre-computation of the smallest $k$ left singular vectors of $D$ to high accuracy.

Better algorithms are typically known for eigenvalue problems, compared to singular value problems. The Dirac operator being $\Gamma_5$-Hermitian [13], i.e. $(\Gamma_5 D)^H = \Gamma_5 D$, we can extract the singular vectors of $D$ from the eigenvectors of $Q = \Gamma_5 D$

$$Q = X\Lambda X^H \Rightarrow D = (\Gamma_5 X\mathrm{sign}(\Lambda))\mathrm{abs}(\Lambda)X^H, \tag{7}$$

with $X$ the eigenvectors of $Q$, $\Lambda$ the diagonal matrix containing the eigenvalues of $Q$, $U = \Gamma_5 X\mathrm{sign}(\Lambda)$ the left singular vectors of $D$ and $V = X$ its right singular vectors. The matrix $\mathrm{abs}(\Lambda)$, which corresponds to the singular values of $D$, is simply the element-wise absolute value of $\Lambda$, and $\mathrm{sign}(\Lambda)$ is such that $\Lambda = \mathrm{sign}(\Lambda)\mathrm{abs}(\Lambda)$.

## 3. Multigrid Multilevel Monte Carlo

The variance reduction technique introduced in [6] makes use of a multigrid hierarchy. We explain now what a multigrid hierarchy is, its relation to inexact deflation, and its use in multigrid multilevel Monte Carlo.

### 3.1 Multigrid

Already with the basic Hutchinson estimator, with or without deflation, LQCD simulations need to solve linear systems of equations of the form $Dx = b$. The Dirac operator $D$ is typically ill-conditioned, with a spectrum that makes it hard for traditional methods such as Krylov-based iterative solvers to find the solution $x$ with tolerable effort. In current LQCD simulations, multigrid solvers [14, 15] are the state of the art when dealing with $Dx = b$ [16–19].

In a two-level multigrid solver, a *smoother* is applied at the original grid (i.e. the finest level), followed by a *coarse-grid correction*. The smoother, which typically consists of a few iterations

of a method such as Gauss-Seidel, Jacobi or GMRES, removes high-frequency components of the error, and the coarse-grid correction serves as a complement to the smoother, in charge of dealing with those modes of the error that have not been dealt with by the smoother. In a more than two-level multigrid method, the coarse grid correction is obtained by applying the two-level method recursively. The operator at level $\ell$ of the resulting multigrid hierarchy is labeled as $D_\ell$. There is an operator that allows the transfer of data from level $\ell$ to $\ell + 1$, known as the restriction operator $R_\ell$; moving data in the opposite direction is done via the interpolation operator $P_\ell$.

Mathematically, one iteration of two-level multigrid can be described by the sequence

$$
\begin{aligned}
r &\leftarrow b - D_1 x \\
x &\leftarrow x + S_1^{(\nu_1)} r \quad \text{(pre-smoothing, } \nu_1 \text{ iterations)} \\
r &\leftarrow b - D_1 x \\
x &\leftarrow x + P_1 D_2^{-1} R_1 r \quad \text{(coarse-grid correction)} \\
r &\leftarrow b - D_1 x \\
x &\leftarrow x + S_1^{(\nu_2)} r \quad \text{(post-smoothing, } \nu_2 \text{ iterations)}
\end{aligned}
\tag{8}
$$

where $r$ is the *residual*, and the operator $S_\ell^{(\nu)}$ is the smoother at level $\ell$, with $\nu$ the number of iterations of the smoother. There are different ways of constructing the coarse-grid operator $D_2$, e.g. the Petrov-Galerkin approach $D_2 = R_1 D_1 P_1$.

If the matrix $D_2$ is still too large, the application $D_2^{-1}(R_1 r)$ in eq. (8) can be further computed via another two-level method. This can be put up in a recursive manner, leading to a multilevel multigrid solver.

The interpolator $P_\ell$ can be built in a geometric or algebraic way. In the former, $P_\ell$ is based on the geometry of the lattice and information of an underlying infinite-dimensional operator prior to discretization . In the latter, $P_\ell$ is rather constructed via the information contained in the matrix $D_\ell$. Due to the random nature of the gauge links in LQCD, algebraic multigrid is necessary to have an effective linear solver. Furthermore, the aggregation-based construction of $P_\ell$ from $D_\ell$ in lattice QCD relies on a concept known as *local coherence* [20], which states that many low modes of $D_\ell$ can be approximately obtained from just a few low modes of the same operator, by looking at the local behaviour of those few modes. The construction of the multigrid hierarchy that we follow here is the one presented in [21].

### 3.2 Inexact Deflation and Multigrid Multilevel Monte Carlo

The orthogonal projector in eq. (4) is built using exact singular vectors $u_i$. If for computational efficiency these are only computed approximately, the computation of $\text{tr}(D^{-1}\Pi) = \text{tr}(D^{-1}UU^*) = \text{tr}(U^*D^{-1}U)$ requires $k$ extra inversions in eq. (5). A cheaper alternative is to use inexact deflation as in [22]. For this, the projector

$$
\Pi = U_k (V_k^H D U_k)^{-1} V_k^H D
\tag{9}
$$

is used. This *oblique* projector splits the trace, using $D^{-1} = (I - \Pi)D^{-1} + \Pi D^{-1}$, as

$$
\text{tr}(D^{-1}) = \text{tr}(D^{-1} - U_k (V_k^H D U_k)^{-1} V_k^H) + \text{tr}((V_k^H D U_k)^{-1} V_k^H U_k)
\tag{10}
$$

The second term in eq. (10) now requires only $k$ matrix-vector multiplications with $D$ acting on the $k$ deflation vectors $u_i$ and the inversion of the small $k \times k$ matrix $V_k^H D U_k$, but no solves with the large matrix $D$. The first term in eq. (10) is expected to have a reduced variance, similar to the case of the orthogonal projector (4) in exact deflation.

In the multigrid context we know from local coherence [20], that the range of the projection $P_1$ is composed of many approximate low modes of $D_1$. This fact is used in [22] to construct inexact deflation with the projector $\Pi_1 = P_1(P_1^H D P_1)^{-1} P_1^H D = P_1 D_2^{-1} P_1^H D$.

Based on the general multilevel Monte Carlo approach [23], the multigrid multilevel Monte Carlo method proposed in [6] applies this construction recursively to obtain the splitting (note that $P_\ell^H P_\ell = I$)

$$\text{tr}(D^{-1}) = \sum_{\ell=1}^{L-1} \text{tr}\left(D_\ell^{-1} - P_\ell D_{\ell+1}^{-1} P_\ell^H\right) + \text{tr}(D_L^{-1}), \tag{11}$$

where $\ell$ runs over the different levels in the multigrid hierarchy, with $\ell = L$ being the coarsest level, and $D_1 = D$. The goal in eq. (11) is to have a sequence of difference-level operators $M_\ell := D_\ell^{-1} - P_\ell D_{\ell+1}^{-1} P_\ell^H$ with reduced variance when computing their traces, and a last term $\text{tr}(D_L^{-1})$ with large variance but cheap to compute.

The expression in eq. (11) has been used in [6] to compute $\text{tr}(D^{-1})$ in the case of the (1+1)-dimensional Schwinger model, in particular. Results for a $128^2$ lattice are displayed here in fig. 1, where $m$ is the mass parameter of the Dirac operator, `cost` is given in FLOPS, and `eps` is the relative tolerance used in the stopping condition. The multigrid multilevel Monte Carlo method presented
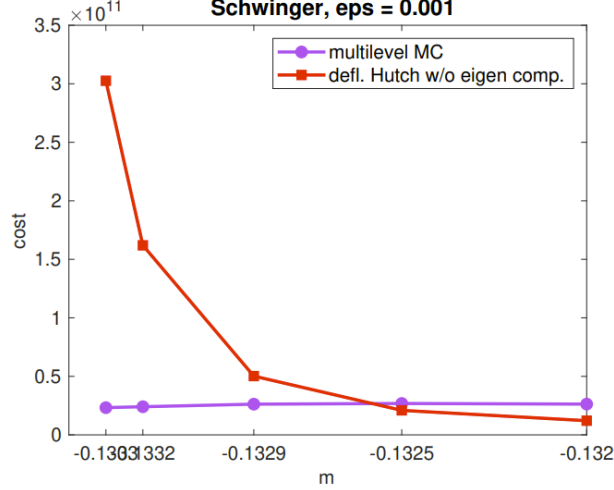


**Figure 1:** Cost versus mass parameter $m$ in the computation of the trace $\text{tr}(D^{-1})$, with `eps` $= 10^{-3}$, for a (1+1)-dimensional Wilson-Schwinger operator $D$ and a lattice of size $128^2$. A four-level multigrid hierarchy is employed. Figure taken from [6].

in [6] displays insensitivity to conditioning, in the example illustrated here in fig. 1, and outperforms a highly tuned exactly deflated Hutchinson by a factor of around 12 for the most ill-conditioned case (i.e. for the smallest value of $m$). Note that we did not include the work for computing the eigenvectors used in exactly deflated Hutchinson in fig. 1.

## 4. Deflated Multigrid Multilevel Monte Carlo

In fig. 1, deflated Hutchinson outperforms multigrid multilevel Monte Carlo for $m = -0.132$ and $m = -0.1325$. We can revert that via two improvements:

- skipping a level

- applying exact deflation on every difference level in multigrid multilevel Monte Carlo

### 4.1 Skipping a level

From the numerical experiments performed in [6] for the Schwinger model, the number of nonzero elements in $D_1$ is approximately the same as in $D_2$. This renders the second level almost as expensive as the first one, which we take into account in this section.

Assuming a four-level multigrid hierarchy, and skipping the second level in the trace decomposition in eq. (11), we can then write

$$\mathrm{tr}(D^{-1}) = \mathrm{tr}(D_1^{-1} - P_1 P_2 A_3^{-1} P_2^H P_1^H) + \mathrm{tr}(A_3^{-1} - P_3 A_4^{-1} P_3^H) + \mathrm{tr}(A_4^{-1}) \tag{12}$$

The first term in eq. (12) might see an increased variance with respect to the first or second terms in the summation in eq. (11), but we might see a gain in total cost due to avoiding inversions with $D_2$ in the multilevel trace expansion.

### 4.2 Exact deflation on difference levels

Following the discussion in section 2, we can apply a similar exact deflation approach on each difference level in multigrid multilevel Monte Carlo. For this, we compute the largest singular vectors for the difference-level operator

$$M_\ell := D_\ell^{-1} - P_\ell D_{\ell+1}^{-1} P_\ell^H. \tag{13}$$

We can reduce the computation of singular triplets to eigenpairs in a way to what we outlined in section 2. We can construct Hermitian difference-level operators using the relations

$$\Gamma_5^\ell P_\ell = P_\ell \Gamma_5^{\ell+1}, \quad P_\ell^H \Gamma_5^\ell = \Gamma_5^{\ell+1} P_\ell^H, \quad (\Gamma_5^\ell)^H = \Gamma_5^\ell, \quad (\Gamma_5^\ell)^H \Gamma_5^\ell = I, \tag{14}$$

which come from the "spin-preserving" algebraic multigrid construction discussed in [17] and which imply

$$Q_\ell := \Gamma_5^\ell D_\ell = D_\ell^H \Gamma_5^\ell = Q_\ell^H.$$

With these relations at hand, we obtain a Hermitian operator which we can use for the indirect extraction of the singular vectors of $M_\ell$:

$$J_\ell = M_\ell \Gamma_5^\ell = D_\ell^{-1} \Gamma_5^\ell - P_\ell D_{\ell+1}^{-1} P_\ell^H \Gamma_5^\ell = Q_\ell^{-1} - P_\ell D_{\ell+1}^{-1} \Gamma_5^{\ell+1} P_\ell^H = Q_\ell^{-1} - P_\ell Q_{\ell+1}^{-1} P_\ell^H. \tag{15}$$

Then, in the same way as in section 2, we extract singular vectors of $M_\ell$ via eigenvectors of $J_\ell$

$$J_\ell = X \Lambda X^H \Rightarrow M_\ell = X \Lambda X^H \Gamma_5^\ell \Rightarrow U = X\mathrm{sign}(\Lambda), S = \mathrm{abs}(\Lambda), V = \Gamma_5 X \tag{16}$$

## 4.3  Results

We have implemented the computation of $\mathrm{tr}(D^{-1})$ via deflated multigrid multilevel Monte Carlo in Python[1]. We have performed a numerical test with the exact same Schwinger matrix used in [6]. We take $m_0 = -0.1320$ here, and we run on a single core of a node with 44 cores Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz. We seek the trace with a relative tolerance of $10^{-4}$, with a relative residual norm of $10^{-12}$ for the linear solves on each sample. When using exact deflation, the relative tolerance in the eigensolver is $10^{-9}$, and $10^{-1}$ for inexact deflation. When calling the eigensolver on $J_\ell$, see eq. (16), the linear solves have a relative tolerance of $10^{-12}$ in exact deflation while they stop at $10^{-3}$ when computing vectors for inexact deflation. The results are presented in table 1.

| Method | Deflation type | Nr. defl. vectors | eig.+direct | trace | total |
|--------|----------------|-------------------|-------------|-------|-------|
| MGMLMC | none | - | - | 276,826.66 | 276,826.66 |
| MGMLMC+skip | none | - | - | 130,208.2 | 130,208.2 |
| Hutchinson | exact | 1,024 | 595.51 | 102,838.06 | 103,433.57 |
| | | 2,048 | 3,337.336 | 102,060.58 | 105,397.916 |
| | | 4,096 | 11,597.91 | 78,727.95 | 90,325.86 |
| | | 8,192 | 77,765.96 | 74,980.76 | 152,746.72 |
| MGMLMC+skip | exact | 512 & 510 | 3,633.41 | 9,903.29 | 13,536.70 |
| | | 1,024 & 510 | 5,878.38 | 7,619.27 | 13,497.65 |
| | | 2,048 & 510 | 13,731.51 | 7,359.11 | 21,090.62 |
| MGMLMC+skip | inexact | 512 & 510 | 2,274.88 | 10,073.90 | 12,348.78 |
| | | 1,024 & 510 | 4,378.46 | 7,588.69 | 11,967.15 |
| | | 2,048 & 510 | 9,959.23 | 7,266.92 | 17,226.15 |

**Table 1:** Execution times (in seconds) when computing $\mathrm{tr}(D^{-1})$ via various methods described in earlier sections. In the third column, which displays the number of deflation vectors used, MGMLMC+skip has only two difference levels due to skipping the second one. The fourth column presents the times to obtain the deflation vectors and in the case of inexact deflation it includes the extra time due to inversions (see the paragraph right before eq. (9) on the need for these inversions), while the fifth column corresponds to the times for computing $\mathrm{tr}(D^{-1})$. The last column shows the total execution time i.e. for eigensolving plus computation of $\mathrm{tr}(D^{-1})$.

Instead of a cost model based on FLOPS, as in [6], we have opted for time measurements here. To reduce overheads due to the interpreter in Python, we measure execution times of very specific operations: $D_\ell$, deflations, $P_\ell$, $P_\ell^H$ and axpy operations.

As has already been reported in [6] and displayed here in fig. 1, exactly deflated Hutchinson clearly outperforms MGMLMC when $m_0 = -0.132$. We revert this behaviour here via deflation on MGMLMC, as illustrated in table 1. In that table, MGMLMC+skip represents skipping the second difference level, which is of great computational benefit because in the multigrid hierarchy for the matrix used here, the number of nonzero elements in $D_1$ and $D_2$ are roughly the same. Skipping the second difference level impacts negatively the variance, and the sample size needed for convergence on the difference level going from $\ell = 1$ to $\ell = 3$ increases with respect to the original two levels

---

[1]The code can be found here.

7

when no skipping is done, but this increase in the sample size is small enough that we still see a large gain when skipping a level, as can be seen from the first two rows in table 1.

Skipping a level is not enough for MGMLMC to outperform exactly deflated Hutchinson, hence we resort to deflation on MGMLMC. As stated before, when eigensolving for exact deflation in the MGMLMC case, i.e. when computing eigenpairs of the operator $J_\ell$ in eq. (16), we use a tolerance of $10^{-9}$ for the eigensolver and of $10^{-12}$ for the solves appearing in $M_\ell$ in each call of $J_\ell$ which happens at each iteration of the eigensolver. We can see in table 1 that, for exactly deflated MGMLMC, the growth in execution time stops being linear at some point, and this is due to asking for a tolerance $10^{-9}$ in the eigensolver, which leads to the orthogonalizations within the eigensolver starting to dominate. This is less pronounced for inexactly deflated MGMLMC, where we ask for a tolerance of $10^{-1}$ from the eigensolver.

This dominance of the orthogonalizations in the eigensolver is expected when a very large number of deflation vectors is set, and we can see it already in exactly deflated Hutchinson in table 1. The size of the L3 cache in the machine where we performed our numerical experiments is 55 MB. A matrix of size $32,768 \times 128$ in double precision, which is the case of the matrix containing the deflation vectors when we deflate 128 of them, has a size of 64 MB, and this is too large to sustain coherence with the largest cache level. When using 128 deflation vectors, though, the multigrid solves are still considerably more expensive than the projections associated to deflation, especially because the Python installation that we use has been compiled with BLAS enabled, and we get BLAS3 performance for these projections.

These projections affect not only the eigensolver, but also the computation of the trace. In table 1, in the case of exactly deflated Hutchinson, we see practically no gain when going from 1,024 to 2,048 deflation vectors, stemming from the deflations going from 11.5% to 18.7% of the overall execution time of the trace, respectively. The same happens in deflated MGMLMC when going from 1,024 & 510 to 2,048 & 510, regardless of the type of deflation.

The best result is, as expected, MGMLMC with the combined use of skipping a level difference plus inexact deflation. With 1,024 & 510 deflation vectors, this method outperforms the best case of exactly deflated Hutchinson by a factor of 7.5 in the overal execution time. This gain comes from the need of less deflation vectors for reducing the variance of each level difference, this coming in turn from the inexact deflation already performed by those difference levels. Furthermore, having less deflation vectors benefits the performance of the method computationally as well, as we have previuosly described.

## 5. Outlook

A first possible improvement for deflated multigrid multilevel Monte Carlo is to switch to an eigensolver more in accordance with the eigenproblem at hand. As described in section 4.2, the computation of the singular vectors is done via eigensolving with the operator $J_\ell = Q_\ell^{-1} - P_\ell Q_{\ell+1}^{-1} P_\ell^H$. This could be converted into a generalized eigenvalue problem with a non-Hermitian operator, which can then be treated via e.g. Jacobi-Davidson in a perhaps more efficient manner. We are currently working on implementing and testing inexactly deflated MGMLMC in the context of lattice QCD.

# References

[1] Arjun Singh Gambhir. *Disconnected Diagrams in Lattice QCD*. PhD thesis, College of William and Mary, 2017.

[2] Andreas Stathopoulos, Jesse Laeuchli, and Kostas Orginos. Hierarchical probing for estimating the trace of the matrix inverse on toroidal lattices. *SIAM J. Sci. Comput.*, 35(5):299–322, 2013.

[3] Jesse Laeuchli and Andreas Stathopoulos. Extending hierarchical probing for computing the trace of matrix inverses. *SIAM J. Sci. Comput.*, 42(3):A1459–A1485, 2020.

[4] M. F. Hutchinson. A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines. *Comm. Statist. Simulation Comput.*, 19(2):433–450, 1990. ISSN 0361-0918. doi: 10.1080/03610919008812864.

[5] Arjun Singh Gambhir, Andreas Stathopoulos, and Kostas Orginos. Deflation as a method of variance reduction for estimating the trace of a matrix inverse. *SIAM J. on Sci. Comput.*, 39 (2):A532–A558, 2017. doi: 10.1137/16M1066361.

[6] Andreas Frommer, Mostafa Nasr Khalil, and Gustavo Ramirez-Hidalgo. A multilevel approach to variance reduction in the stochastic estimation of the trace of a matrix. *SIAM Journal on Scientific Computing*, 44(4):A2536–A2556, 2022.

[7] S.J. Dong and K.F. Liu. Stochastic estimation with $Z_2$ noise. *Phys. Lett. B*, 328:130–136, 1994.

[8] Walter Wilcox. Noise methods for flavor singlet quantities. In *Numerical Challenges in Lattice Quantum Chromodynamics*, pages 127–141. Springer, 2000.

[9] Constantia Alexandrou, S Bacchio, M Constantinou, J Finkenrath, K Hadjiyiannakou, K Jansen, G Koutsou, and A Vaquero Aviles-Casco. Proton and neutron electromagnetic form factors from lattice qcd. *Physical Review D*, 100(1):014509, 2019.

[10] C Alexandrou, S Bacchio, M Constantinou, J Finkenrath, K Hadjiyiannakou, K Jansen, G Koutsou, H Panagopoulos, G Spanoudes, Extended Twisted Mass Collaboration, et al. Complete flavor decomposition of the spin and momentum fraction of the proton using lattice qcd simulations at physical pion mass. *Physical Review D*, 101(9):094513, 2020.

[11] Zohreh Davoudi, William Detmold, Phiala Shanahan, Kostas Orginos, Assumpta Parreno, Martin J Savage, and Michael L Wagman. Nuclear matrix elements from lattice qcd for electroweak and beyond-standard-model processes. *Physics Reports*, 900:1–74, 2021.

[12] Lloyd N Trefethen and David Bau III. *Numerical Linear Algebra*, volume 50. SIAM, 1997.

[13] C. Gattringer and C. B. Lang. *Quantum Chromodynamics on the Lattice*, volume 788 of *Lect. Notes Phys.* Springer, 2009.

[14] John W Ruge and Klaus Stüben. Algebraic multigrid. In *Multigrid methods*, pages 73–130. SIAM, 1987.

[15] William L Briggs, Van Emden Henson, and Steve F McCormick. *A Multigrid Tutorial*. SIAM, 2000.

[16] J. C. Osborn, R. Babich, J. Brannick, R. C. Brower, M. A. Clark, S. D. Cohen, and C. Rebbi. Multigrid solver for clover fermions. *PoS*, LATTICE2010:037, 2010. doi: 10.22323/1.105. 0037.

[17] Andreas Frommer, Karsten Kahl, Stefan Krieg, Björn Leder, and Matthias Rottmann. Adaptive aggregation-based domain decomposition multigrid for the lattice Wilson-Dirac operator. *SIAM journal on scientific computing*, 36(4):A1581–A1608, 2014.

[18] Richard C Brower, K Moriarty, E Myers, and Claudio Rebbi. The multigrid method for fermion calculations in quantum chromodynamics. *Multigrid Methods: Theory, Applications, and Supercomputing, SF McCormick, ed*, 110:85–100, 1987.

[19] A. Frommer, K. Kahl, S. Krieg, B. Leder, and M. Rottmann. Adaptive aggregation based domain decomposition multigrid for the lattice Wilson-Dirac operator. *SIAM J. Sci. Comp.*, 36(4):A1581–A1608, 2014.

[20] Martin Lüscher. Local coherence and deflation of the low quark modes in lattice QCD. *Journal of High Energy Physics*, 2007(07):081, 2007.

[21] Matthias Rottmann. *Adaptive Domain Decomposition Multigrid for Lattice QCD*. PhD thesis, Wuppertal U., 2016.

[22] Eloy Romero, Andreas Stathopoulos, and Kostas Orginos. Multigrid deflation for lattice QCD. *Journal of Computational Physics*, 409:109356, 2020.

[23] Michael B Giles. Multilevel Monte Carlo methods. *Acta Numer.*, 24:259–328, 2015. doi: 10.1017/S096249291500001X.