# Quasi-Monte Carlo methods for computing flow in random porous media

*I. G. Graham, F. Y. Kuo, D. Nuyens, R. Scheichl, and I. H. Sloan*

# Bath Institute For Complex Systems

# Quasi-Monte Carlo methods for computing flow
# in random porous media

I. G. Graham[1,4], F. Y. Kuo[2], D. Nuyens[2,3], R. Scheichl[1], and I. H. Sloan[2]

[1] Dept of Mathematical Sciences, University of Bath, Bath BA2 7AY UK
`I.G.Graham@bath.ac.uk`, `R.Scheichl@bath.ac.uk`

[2] School of Mathematics and Statistics, University of NSW, Sydney NSW 2052, Australia.
`f.kuo@unsw.edu.au`, `i.sloan@unsw.edu.au`

[3] Department of Computer Science, K.U.Leuven, Celestijnenlaan 200A, B-3001 Heverlee.
`dirk.nuyens@cs.kuleuven.be`

[4] Corresponding author: Telephone: +44, 1225 386989, Fax: +44, 1225 386492

## Abstract

We devise and implement quasi-Monte Carlo methods for computing the expectations of nonlinear functionals of random fields arising in the modeling of fluid flow in random porous media. Specific examples include the effective permeability of a block of rock, the pressure head at a chosen point and the breakthrough time of a pollution plume being convected by the velocity field. The mathematical model is a system of first order partial differential equations in space with a random field describing the permeability. Our emphasis is on situations where a very large number of dimensions are necessary to obtain a reasonable accuracy in probability space, and where classical Monte Carlo methods with random sampling are currently the method of choice. As an alternative we introduce quasi-Monte Carlo methods which use deterministically chosen sample points in probability space. Our algorithm performs finite element approximation for each realization of the field, and approximates the necessary element integrals by using values of the field computed only on a suitable regular grid, with the help of FFT techniques. In this way we avoid the use of a truncated Karhunen-Loève expansion, but introduce high nominal dimension in probability space. Numerical experiments with 2-dimensional rough random fields, high variance and small length scale are reported, showing that quasi-Monte Carlo method consistently outperforms the Monte Carlo method, with a noticeably better than $\mathcal{O}(N^{-1/2})$ convergence rate and a smaller implied constant, where $N$ is the number of samples. Moreover, the rate of convergence of quasi-Monte Carlo methods does not appear to degrade as the nominal dimension increases. Examples with dimension as high as $10^6$ are reported.

**Keywords:** Quasi-Monte Carlo, High Dimensional Quadrature, Fluid Flow, Random Porous Media, Circulant Embedding, Fast Fourier Transform

## 1 Introduction

Many physical, biological or geological models involve spatially varying input data which may be subject to uncertainty. This induces a corresponding uncertainty in the outputs of the model and in any physical quantities of interest which may be derived from these outputs. A common way to deal with these uncertainties is by considering the input data to be a *random field*, in which case the derived quantity of interest is either a random variable (e.g., the effective conductivity of a composite material) or another random field (e.g., the velocity of fluid in a random medium). The computational goal is usually to find the expected value, higher order moments or other statistics of these derived quantities.

The aim of this paper is to formulate and implement *quasi-Monte Carlo* (QMC) methods for the computation of flow in random media. There is a huge interest in problems of this sort in the geosciences, where *Monte Carlo* (MC) methods are regularly employed, see, e.g.,

[10, 25, 26, 44]. We shall assume (as often done in practice) that the permeability is a lognormal Gaussian random field, and, as examples of quantities of physical interest, we shall study the pressure head at a certain point, the effective permeability, and the breakthrough time of a plume of pollution moving in the fluid.

The expected value or higher moments of these physical quantities are, by definition, integrals. These are in principle infinite dimensional, but in practice they are approximated by finite (but often very high dimensional) integrals. MC methods for the flow problem approximate such integrals by equal-weight quadrature rules with randomly chosen points. QMC methods are also equal-weight quadrature rules but aim to outperform the MC methods by clever choices of deterministic points. As we shall show in our numerical experiments, our QMC methods consistently outperform MC methods in all cases studied, and in the best cases enjoy an improvement in running time by a factor of several hundred. Moreover, the QMC methods are observed to converge noticeably faster than the rate $\mathcal{O}(N^{-1/2})$, where $N$ is the number of samples, and this rate is robust as the dimension $d$ of the domain of integration increases. (Some of our computations are even for dimension $d \approx 10^6$.) By contrast MC methods are also robust to increasing dimension but converge only with the well-known characteristic rate of $\mathcal{O}(N^{-1/2})$. In our theoretical discussion in §4.1 we review sufficient conditions for the observed dimension-independent convergence behavior of QMC methods. However it is not known if these sufficient conditions are satisfied by the complicated flow problems solved here and so this aspect of the presentation is experimental.

In recent years there has been a great interest in methods that treat PDE problems with random input data by simultaneous approximation in both physical and probability space. These go under names such as *stochastic Galerkin*, *stochastic collocation* or *polynomial chaos*, see, e.g., [3, 2, 15, 29, 30, 33, 37]. Many of these methods start by representing the random field using the *Karhunen-Loève* (KL) expansion, and then truncate this expansion after a finite number of terms before discretization in space. From another point of view, it is sometimes said that these methods depend on a "finite-dimensional noise" assumption. While such approaches can be very effective when the KL expansion converges rapidly, they face the serious challenges of high cost combined with large truncation error when the convergence of the KL expansion is slow. Such slow convergence is a common feature in practical flow problems and motivates the development of the alternative methods proposed in this paper.

The starting point of our approach is the observation that if we first approximate the random PDE problem ((2.1)–(2.3) below) by finite elements in space then, without paying a penalty in the order of finite element accuracy, the resulting (still random, but now) discrete problem may be formulated using data taken from only the *finite* set of values of the original continuous field sampled on a regular grid. Thus realizations of the random discrete problem may be computed by realizing only the discrete field, whose covariance *matrix* is inherited from the covariance function of the continuous field.

While this observation is hardly deep, it turns out to be crucially important to the design of efficient algorithms. Fast linear algebra techniques can be designed to sample random vectors with a given covariance matrix, thus avoiding the truncations required in a KL-based approach. In fact these linear algebra techniques can be based on existing methods for realizing random fields at discrete points in space. In our implementation we use the circulant embedding approach combined with FFT for homogeneous Gaussian random fields, see [5, 11]. It produces realizations of such fields at uniformly spaced points on a tensor product grid. All our experiments here are on $(0, 1)^2$, but the method is applicable to more general domains, as long as the grid on which the random field is sampled is a regular tensor product grid.

In our approach the dimension $d$ of the domain of the integral being computed is linked to the dimension $M$ of the finite element space, in fact $d = \mathcal{O}(M)$. In practice this relatively large dimension does not seem to have any detrimental effect on the performance of the QMC methods. The (extended) covariance matrix is diagonalized (using FFT), enabling an ordering

of the variables in decreasing importance which is crucial to the success of QMC. Each point evaluation of the integrand requires (i) the computation of an $M$-dimensional vector representing the permeability field at the finite element grid points through diagonalizing a $d \times d$ circulant matrix which (via FFT) costs $\mathcal{O}(d \log d) = \mathcal{O}(M \log M)$ operations and (ii) the solution of the mixed finite element (FE) approximation of the relevant system of PDEs, which results in an indefinite highly ill-conditioned system which we solve robustly in $\mathcal{O}(M)$ operations. Thus the overall computational cost is $\mathcal{O}(NM \log M)$, where $N$ is the number of quadrature points (i.e., the number of times the random field is sampled). We choose mixed FEs because they provide physically realistic flow fields which are crucial to the computation of quantities of physical interest. We combine this discretization with the divergence-free reduction technique from [7, 35] which provides a robust solver and fast computation of all quantities of physical interest.

The plan of the paper is as follows. In §2 we outline the target problems arising from flow in random porous media. In §3 we describe the linear algebra approach to generating the random discrete field and show how this leads to formulas for expectations of quantities of interest in terms of integrals over a (possibly high dimensional) unit cube. We also explain how the QMC (and MC) methods can be used to evaluate these high-dimensional integrals and mention here how our method is related to methods based on truncating the KL expansion, such as the stochastic collocation method. In §4 we provide more details about QMC methods and give an overview of the expected performance and strategies for implementing these methods. The following two sections are devoted to fast implementation: §5 explains how circulant embedding techniques can be used to efficiently compute the formulas developed in §3, while §6 discusses fast methods for computing quantities of interest using mixed finite elements. In §6 we restrict the physical domain to a simple rectangular "flow cell" and study the pressure head at a point, the effective permeability and the breakthrough time as quantities of physical interest. Extensive numerical experiments on the flow cell are given in §7 which confirm the efficiency and robustness of our QMC method. For example, we solve a 2-dimensional problem (with parameters as in Case 5 in Table 1 of §7) in which the contributions from the first 300 terms in the KL expansion are all greater than one tenth of the largest term (i.e., $\sqrt{\mu_{300}} > 0.1\sqrt{\mu_1}$ in the KL expansion (3.10) below). Nevertheless, we are able to compute the effective permeability for this problem to an error of $10^{-3}$ with respect to both physical space and probability space by avoiding a truncation of the KL expansion. To achieve this accuracy we require a spatial mesh with $M \approx 10^6$ degrees of freedom and about $N \approx 2.9 \times 10^4$ samples, and the computation takes about 35 hours on a standard laptop computer. In comparison, to obtain the same accuracy with MC, we require about $N \approx 5.3 \times 10^5$ samples or more than 2 months on a standard laptop. Note that CPU times are of course greatly reduced on multi-processor machines (which we used), due to the independence of the PDE problems for individual realizations of the random field and the resulting optimal parallel scalability of our method.

## 2    Random PDE model and its spatial discretization

We consider flow in a 2D porous medium governed by *Darcy's law*:

$$\vec{q} + k\,\vec{\nabla}p = \vec{0}, \tag{2.1}$$

and the law of mass conservation:

$$\vec{\nabla} \cdot \vec{q} = 0, \tag{2.2}$$

where $\vec{q}, p$ are respectively the (2-dimensional) velocity (also called the specific discharge) and (scalar) residual pressure, which are to be found on a bounded domain $D \subset \mathbb{R}^2$. The boundary conditions are taken to be

$$p = g \quad \text{on} \quad \Gamma^D, \qquad \vec{q} \cdot \vec{n} = 0 \quad \text{on} \quad \Gamma^N \tag{2.3}$$

where $\Gamma^D$ and $\Gamma^N$ are the Dirichlet and Neumann parts of $\Gamma$, the boundary of $D$, and $\vec{n}$ denotes the outer unit normal on $\Gamma$. Note that although we restrict ourselves here to 2 space dimensions, only the fast mixed finite element solver in §6 is specific to 2D. By replacing this with a 3D mixed finite element solver (many exist, but not as fast), our approach will extend in a straightforward way to 3D.

## 2.1 Modeling the permeability as a random field

In (2.1) $k$ is the permeability (more precisely the ratio of permeability to dynamic viscosity) and is modeled as a random field. Because $k$ is physically positive, we shall here make the popular and natural assumption that $k$ is a *lognormal* random field, i.e.,

$$k(\vec{x}\,;\,\omega) \;=\; \exp(Z(\vec{x}\,;\,\omega))\,, \tag{2.4}$$

where $Z = Z(\vec{x}\,;\,\cdot)$, is a zero mean Gaussian random field on $D$, with a specified continuous covariance function $r(\vec{x}, \vec{y})$, and $\omega$ denotes an event in the probability space $(\Omega, \mathcal{F}, \mathcal{P})$. That is

$$\mathbb{E}(Z(\vec{x}\,;\,\cdot)) \;=\; 0 \quad \text{and} \quad \mathbb{E}(Z(\vec{x}\,;\,\cdot)\,Z(\vec{y}\,;\,\cdot)) \;=\; r(\vec{x}, \vec{y}) \quad \text{for all} \quad \vec{x}, \vec{y} \in D,$$

where $\mathbb{E}$ denotes expectation with respect to the Gaussian probability measure on $\Omega$.

There is some evidence from field data that (2.4) gives a reasonable representation of reality in certain cases (see [13, 20]). Note that some authors add a positive constant to the right-hand side of (2.4), which ensures that the resulting PDE for the pressure is uniformly coercive in probability space. This turns out not to be necessary in practice and we avoid it here.

Throughout we will assume that $Z$ is *homogeneous* (see, e.g., [1, p.24]), i.e., its covariance function satisfies

$$r(\vec{x}, \vec{y}) \;=\; \rho(\vec{x} - \vec{y}),$$

where $\rho : \mathbb{R}^2 \to \mathbb{R}$ is a suitably behaved given function. An example of particular significance is

$$\rho(\vec{t}) \;=\; \sigma^2 \, \exp\Big( -\,\|\vec{t}\|_p/\lambda \Big), \tag{2.5}$$

with $\|\cdot\|_p$ denoting the $\ell_p$ norm on $\mathbb{R}^2$, and where the parameters $\sigma^2$ and $\lambda$ denote respectively the *variance* and (*correlation*) *length scale*. All our computations will be for this covariance function with $p = 1$ or $p = 2$. Since the function $\rho$ is not smooth at the origin, realizations of $Z$ can be quite irregular. In fact Kolmogorov's theorem [1, Theorem 8.3.2] implies that with probability 1, realizations $Z(\vec{x}\,;\,\omega)$ are Hölder continuous with respect to $\vec{x}$, with Hölder exponent $\alpha$ limited to lie in the range $\alpha \in [0, 1/2)$. Decreasing the length scale increases the frequency of oscillations in $Z$, while increasing the variance increases their amplitude. The roughness in $Z$ induces roughness in $k$, and in the solution $(\vec{q}, p)$, and, as we shall see, increases the number of degrees of freedom needed to achieve acceptable accuracy in the solution of (2.1)–(2.3).

Our overall goal is to compute expectations of a random variable $\mathcal{G}(Z)$, derived from $Z$ through solving the PDE problem (2.1)–(2.3) with $k$ given by (2.4). A simple example of physical interest is the pressure head at a particular point $\vec{x}^* \in D$, i.e.,

$$\mathcal{G}(Z) \;=\; p(\vec{x}^*)\,, \tag{2.6}$$

where $p$ is the solution of (2.1)–(2.3). The computation of each realization $\mathcal{G}(Z)$ requires the construction of a realization of $Z$, the solution of the PDE system (2.1)–(2.3), with permeability given by (2.4), and finally the evaluation of the pressure approximation at $\vec{x}^*$. We will give other examples of $\mathcal{G}(Z)$ in §6.

## 2.2 Solving the PDE problem using the mixed finite element method

To compute approximate solutions to (2.1)–(2.3) for a given realization of (2.4), we use the mixed finite element method. Problem (2.1)–(2.3) is written in weak form as the problem of seeking $(\vec{q}, p) \in H_{0,N}(\mathrm{div}, D) \times L_2(D)$,

$$
\left.
\begin{array}{rclcll}
m_\omega(\vec{q}, \vec{v}) & + & b(p, \vec{v}) & = & G(\vec{v}), & \text{for all } \vec{v} \in H_{0,N}(\mathrm{div}, D), \\
b(w, \vec{q}) & & & = & 0, & \text{for all } w \in L_2(D),
\end{array}
\right\} \tag{2.7}
$$

where $H_{0,N}(\mathrm{div}, D) := \{\vec{v} \in L_2(D)^2 : \vec{\nabla} \cdot \vec{v} \in L_2(D) \text{ and } \vec{v} \cdot \vec{n} = 0 \text{ on } \Gamma^N\}$,

$$
m_\omega(\vec{q}, \vec{v}) \;=\; \int_D k(\vec{x}; \omega)^{-1}\, \vec{q}(\vec{x}) \cdot \vec{v}(\vec{x})\, \mathrm{d}\vec{x}, \tag{2.8}
$$

$$
b(p, \vec{v}) \;=\; -\int_D p(\vec{x})\, \vec{\nabla} \cdot \vec{v}(\vec{x})\, \mathrm{d}\vec{x}, \quad \text{and} \quad G(\vec{v}) \;=\; -\int_{\Gamma^D} g(\vec{x})\, \vec{v}(\vec{x}) \cdot \vec{n}(\vec{x})\, \mathrm{d}\Gamma(\vec{x}). \tag{2.9}
$$

The subscript $\omega$ in (2.8) indicates that the randomness appears in this bilinear form. In a proper probabilistic setting we would seek random solutions $(\vec{q}, p)$ in an appropriate tensor product space (e.g., $p(\vec{x}) = p(\vec{x}; \omega)$ would be required to be square integrable with respect to Gaussian probability measure on $\Omega$, with values in $L_2(D)$) and (2.7) would be required to be satisfied almost surely for $\omega \in \Omega$). A precise probabilistic setting of the problem can be found, e.g., in [29], but we suppress this here as it is not needed for the formulation of our algorithm.

To discretize (2.7), we introduce a mesh $\mathcal{T}_h$ (here taken to be triangular) on $D$, and we approximate $\vec{q}(\vec{x})$ by $\vec{q}_h(\vec{x}) \in \mathcal{V}_h$, where $\mathcal{V}_h \subset H_{0,N}(\mathrm{div}, D)$ is the space of lowest order *Raviart-Thomas elements* characterized by: (i) $\vec{q}_h(\vec{x}) = \vec{\alpha}_\tau + \gamma_\tau \vec{x}$, $\vec{x} \in \tau$, for each $\tau \in \mathcal{T}_h$ and for some suitable coefficients $\vec{\alpha}_\tau \in \mathbb{R}^2$ and $\gamma_\tau \in \mathbb{R}$; (ii) $\vec{q}_h$ has continuous normal component across the interior edges of the mesh; and (iii) $\vec{q}_h \cdot \vec{n} = 0$ on $\Gamma^N$. (For more details on Raviart-Thomas elements see [4, 24].) The pressure $p(\vec{x})$ is approximated by $p_h(\vec{x}) \in \mathcal{W}_h$, the subspace of $L_2(D)$ consisting of all piecewise constant functions with respect to $\mathcal{T}_h$. The approximate solution $(\vec{q}_h, p_h)$ is then computed as the solution to the discrete system

$$
\left.
\begin{array}{rclcll}
m_\omega(\vec{q}_h, \vec{v}_h) & + & b(p_h, \vec{v}_h) & = & G(\vec{v}_h), & \text{for all } \vec{v}_h \in \mathcal{V}_h, \\
b(w_h, \vec{q}_h) & & & = & 0, & \text{for all } w_h \in \mathcal{W}_h.
\end{array}
\right\} \tag{2.10}
$$

Again we should more properly write $\vec{q}_h(\vec{x}) = \vec{q}_h(\vec{x}; \omega)$ and $p_h(\vec{x}) = p_h(\vec{x}; \omega)$ and require (2.10) to hold almost surely for $\omega \in \Omega$.

Choosing bases $\{\vec{v}_j : j = 1, \ldots, n_v\}$ and $\{w_\ell : \ell = 1, \ldots, n_w\}$ for $\mathcal{V}_h$ and $\mathcal{W}_h$, respectively, and writing $\vec{q}_h := \sum_{j=1}^{n_v} Q_j \vec{v}_j$ and $p_h := \sum_{\ell=1}^{n_w} P_\ell w_\ell$, we obtain the symmetric indefinite system

$$
\begin{pmatrix} M_\omega & B \\ B^\mathsf{T} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{Q} \\ \mathbf{P} \end{pmatrix} = \begin{pmatrix} \mathbf{g} \\ \mathbf{0} \end{pmatrix} \quad \in \quad \mathbb{R}^{n_v + n_w} \tag{2.11}
$$

with $B_{i,\ell} = b(w_\ell, \vec{v}_i)$ and $g_i = G(\vec{v}_i)$. The random matrix $M_\omega$ is given by

$$
(M_\omega)_{i,j} \;=\; m_\omega(\vec{v}_j, \vec{v}_i) \;=\; \int_D k(\vec{x}; \omega)^{-1}\, \vec{v}_j(\vec{x}) \cdot \vec{v}_i(\vec{x})\, \mathrm{d}\vec{x}, \quad i, j = 1, \ldots, n_v. \tag{2.12}
$$

As we have remarked above, for the covariance function (2.5), with probability 1, we have $Z(\cdot; \omega) \in C^\alpha(D)$ for $\alpha < 1/2$. It follows from the arguments in [8] that then $\vec{q}_h$ converges to $\vec{q}$ $\omega$-almost-surely in $\Omega$, with order $\mathcal{O}(h^\alpha)$ as $h \to 0$ and, moreover, no degradation of this rate of convergence is incurred if the integrals (2.12) are approximated by even a low order quadrature rule using evaluations of the integrand at one or more points in (or near) $\tau$. As a simple example, if the FE grid is itself regular we may approximate $Z(\vec{x}; \omega)$ in each element $\tau \in \mathcal{T}_h$ by the average of its values at the three nodes of $\tau$, without sacrificing accuracy. Thus

if $\vec{x}_1, \ldots, \vec{x}_M$ are the nodes of the mesh $\mathcal{T}_h$ which lie on $D \cup \Gamma^N$, and if $\boldsymbol{Z} \in \mathbb{R}^M$ denotes the random vector

$$\boldsymbol{Z} = (Z_1, \ldots, Z_M)^{\mathsf{T}} := (Z(\vec{x}_1; \omega), \ldots, Z(\vec{x}_M; \omega))^{\mathsf{T}}, \tag{2.13}$$

then an appropriate approximation of $M_\omega$ may be taken as $\widetilde{M}(\boldsymbol{Z})$, with

$$(\widetilde{M}(\boldsymbol{Z}))_{i,j} := \sum_\tau \bar{k}_\tau^{-1} \int_\tau \vec{v}_i(\vec{x}) \cdot \vec{v}_j(\vec{x}) \, d\vec{x}, \quad \text{where} \quad \bar{k}_\tau := \frac{1}{3} \sum_{\ell : \vec{x}_\ell \in \tau} \exp(Z_\ell), \tag{2.14}$$

and it suffices to take the first sum over all $\tau \subset \operatorname{supp}(\vec{v}_i) \cap \operatorname{supp}(\vec{v}_j)$. The resulting approximate saddle point system is

$$\begin{pmatrix} \widetilde{M}(\boldsymbol{Z}) & B \\ B^{\mathsf{T}} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{Q} \\ \mathbf{P} \end{pmatrix} = \begin{pmatrix} \mathbf{g} \\ \mathbf{0} \end{pmatrix} \qquad \in \quad \mathbb{R}^{n_v + n_w} \tag{2.15}$$

Thus the expected value of a quantity of interest $\mathcal{G}(Z)$ can be approximated by the expected value of its finite element approximation $\mathcal{G}_h(\boldsymbol{Z})$, which only depends on the random vector $\boldsymbol{Z}$. In the particular example (2.6), a realization of $\mathcal{G}(Z)$ would be approximated by

$$\mathcal{G}_h(\boldsymbol{Z}) := p_h(\vec{x}^*)$$

with $p_h$ computed from (2.15).

# 3 Sampling $\boldsymbol{Z}$ and evaluating expectations

The vector $\boldsymbol{Z}$ in (2.13) is Gaussian with mean zero and $M \times M$ positive definite covariance matrix

$$R = \mathbb{E}(\boldsymbol{Z}\boldsymbol{Z}^{\mathsf{T}}) = \left( r(\vec{x}_i, \vec{x}_j) \right)_{i,j=1}^M .$$

A procedure for sampling $\boldsymbol{Z}$ can be based on any real factorization of $R$ in the form

$$R = \Theta\Theta^{\mathsf{T}}, \tag{3.1}$$

where $\Theta$ is a real $M \times M$ matrix (e.g., a Cholesky factorization). From this, it can easily be seen that

$$\boldsymbol{Z} := \Theta\boldsymbol{Y} \tag{3.2}$$

defines a suitable realization, provided $\boldsymbol{Y} := (Y_1(\omega), \ldots, Y_M(\omega))^{\mathsf{T}}$ is a vector of independent standard Gaussian random variables. In contrast to other approaches, such as the truncated KL expansion (see §3.3 below), the formula (3.2) represents the field exactly at the discrete points $\vec{x}_1, \ldots, \vec{x}_M$ without any truncation.

## 3.1 MC and QMC approximations

With the finite element approximation $\mathcal{G}_h(\boldsymbol{Z})$ as described in §2.2, the MC (simulation) method approximates the expected value by an equal-weight average

$$\mathbb{E}(\mathcal{G}(Z)) \overset{\text{FE}}{\approx} \mathbb{E}(\mathcal{G}_h(\boldsymbol{Z})) \overset{\text{MC}}{\approx} \frac{1}{N} \sum_{n=1}^N \mathcal{G}_h(\Theta\boldsymbol{y}^{(n)}), \tag{3.3}$$

where $\boldsymbol{y}^{(1)}, \ldots, \boldsymbol{y}^{(N)} \in \mathbb{R}^M$ are $N$ independent samples of standard Gaussian random vectors. The implementation of this requires the solution of $N$ different instances of (2.15) with $\boldsymbol{Z} = \Theta\boldsymbol{y}^{(n)}$, $n = 1, \ldots, N$.

To apply QMC methods we need to step away from the MC simulation point of view and express the expected value explicitly as an $M$-dimensional integral

$$\mathbb{E}(\mathcal{G}_h(\boldsymbol{Z})) \;=\; \int_{\mathbb{R}^M} \mathcal{G}_h(\boldsymbol{z}) \frac{\exp(-\frac{1}{2}\boldsymbol{z}^{\mathsf{T}} R^{-1}\boldsymbol{z})}{(2\pi)^{M/2}(\det R)^{1/2}} \, \mathrm{d}\boldsymbol{z} \,. \tag{3.4}$$

Moreover, since QMC methods are (equal-weight) quadrature rules defined over the unit cube, it is necessary to transform (3.4) to an integral over the $M$-dimensional unit cube. This can easily be done if we have factorized $R$ in the form (3.1). We introduce the univariate standard normal cumulative distribution function

$$\Phi(y) \;:=\; \int_{-\infty}^{y} \frac{\exp\left(-t^2/2\right)}{\sqrt{2\pi}} \, \mathrm{d}t \qquad \text{for} \quad y \in \mathbb{R},$$

and set $\boldsymbol{\Phi}_M^{-1}(\boldsymbol{x}) := (\Phi^{-1}(x_1), \ldots, \Phi^{-1}(x_M))^{\mathsf{T}} \in \mathbb{R}^M$ for $\boldsymbol{x} \in [0,1]^M$. The successive changes of variables $\boldsymbol{z} = \Theta\boldsymbol{y}$ and $\boldsymbol{y} = \boldsymbol{\Phi}_M^{-1}(\boldsymbol{x})$ transform (3.4) into

$$\mathbb{E}(\mathcal{G}_h(\boldsymbol{Z})) \;=\; \int_{[0,1]^M} \mathcal{G}_h\big(\Theta\boldsymbol{\Phi}_M^{-1}(\boldsymbol{x})\big) \, \mathrm{d}\boldsymbol{x} \,.$$

Then a QMC method approximates the expected value by

$$\mathbb{E}(\mathcal{G}(Z)) \;\overset{\text{FE}}{\approx}\; \mathbb{E}(\mathcal{G}_h(\boldsymbol{Z})) \;\overset{\text{QMC}}{\approx}\; \frac{1}{N} \sum_{n=1}^{N} \mathcal{G}_h\big(\Theta\boldsymbol{\Phi}_M^{-1}(\boldsymbol{x}^{(n)})\big), \tag{3.5}$$

where $\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(N)} \in [0,1]^M$ are $N$ deterministically chosen points from the unit cube. We will return to discuss QMC methods in more detail in §4, and in particular, we will replace (3.5) by a randomized QMC method for practical error estimation.

We remark that the MC method can also be viewed as an equal-weight quadrature rule over the unit cube as in (3.5), but with $\{\boldsymbol{x}^{(n)}\}$ being independent random vectors drawn from a *uniform* distribution in $[0,1]^M$. This latter point of view for MC is less convenient than (3.3) in practice, since it requires the evaluation of $\boldsymbol{\Phi}_M^{-1}$ which has no analytic form.

## 3.2  Embedding in a larger matrix

Both MC and QMC approximations discussed above rely on a factorization of $R$ in the form (3.1). However such factorization can be expensive — a Cholesky factorization requires $\mathcal{O}(M^3)$ operations in general — so we seek more efficient methods by first extending $R$ to a larger symmetric positive definite $d \times d$ matrix $C$. The disadvantage of the larger dimension is offset by the extra freedom in choosing the structure of $C$ in such a way as to make it easy to factorize. Dietrich and Newsam [11] and Chan and Wood [5] proposed to construct $C$ as a circulant extension of $R$, for then a factorization of $C$ can be constructed in $\mathcal{O}(d \log d)$ operations using a fast Fourier transform. Furthermore, this yields an orthogonal diagonalization of $C$ and the resulting eigenstructure allows us to identify the order of importance of the variables, something that is crucial to the success of QMC, see §4. We return to circulant embedding in §5.

The justification provided in [11] for the embedding approach is expressed in terms of obtaining the correct correlation structure in the random MC samples. Since QMC methods are quadrature rules rather than simulation techniques, the justification behind the embedding approach has to be quite different in nature. We provide this justification below. Lemma 1 shows that the extension from $R$ to $C$ still yields an integral of an analogous form to (3.4). Lemma 2 then shows how such integrals can be transformed to the $d$-dimensional unit cube. Corollary 3 applies the result (slightly generalized) to (3.4).

**Lemma 1** *Suppose that $R$ is any symmetric positive definite $M \times M$ matrix and suppose $C$ is a $d \times d$ symmetric positive definite matrix of the form*

$$C = \begin{bmatrix} R & U \\ U^{\mathsf{T}} & V \end{bmatrix} . \tag{3.6}$$

*Then for any integrable function $g : \mathbb{R}^M \to \mathbb{R}$ we have*

$$\int_{\mathbb{R}^M} g(\boldsymbol{z}) \frac{\exp(-\frac{1}{2}\boldsymbol{z}^{\mathsf{T}} R^{-1} \boldsymbol{z})}{(2\pi)^{M/2}(\det R)^{1/2}} \, \mathrm{d}\boldsymbol{z} = \int_{\mathbb{R}^d} g(\boldsymbol{u}_{[1:M]}) \frac{\exp(-\frac{1}{2}\boldsymbol{u}^{\mathsf{T}} C^{-1} \boldsymbol{u})}{(2\pi)^{d/2}(\det C)^{1/2}} \, \mathrm{d}\boldsymbol{u} ,$$

*where, for any vector $\boldsymbol{u} \in \mathbb{R}^d$, $\boldsymbol{u}_{[1:M]} \in \mathbb{R}^M$ is the vector containing the first $M$ components of $\boldsymbol{u}$.*

**Proof.** Let $R = LL^{\mathsf{T}}$ be a Cholesky factorization of $R$. Then with the substitution $\boldsymbol{z} = L\boldsymbol{w}$ we have

$$\int_{\mathbb{R}^M} g(\boldsymbol{z}) \frac{\exp(-\frac{1}{2}\boldsymbol{z}^{\mathsf{T}} R^{-1} \boldsymbol{z})}{(2\pi)^{M/2}(\det R)^{1/2}} \, \mathrm{d}\boldsymbol{z} = \int_{\mathbb{R}^M} g(L\boldsymbol{w}) \frac{\exp(-\frac{1}{2}\boldsymbol{w}^{\mathsf{T}} \boldsymbol{w})}{(2\pi)^{M/2}} \, \mathrm{d}\boldsymbol{w} .$$

A Cholesky factorization of $C$ can be written in the form

$$C = \widetilde{L}\,\widetilde{L}^{\mathsf{T}}, \quad \text{where} \quad \widetilde{L} = \begin{bmatrix} L & 0 \\ K & L' \end{bmatrix} .$$

Hence, introducing the new variable $\boldsymbol{v} = [\boldsymbol{w}^{\mathsf{T}}, \boldsymbol{w}'^{\mathsf{T}}]^{\mathsf{T}}$, where $\boldsymbol{w}' \in \mathbb{R}^{d-M}$ is a vector of dummy variables, and then making the substitution $\boldsymbol{v} = \widetilde{L}^{-1}\boldsymbol{u}$ where $\boldsymbol{u} \in \mathbb{R}^d$, we obtain

$$\begin{aligned}
\int_{\mathbb{R}^M} g(\boldsymbol{z}) \frac{\exp(-\frac{1}{2}\boldsymbol{z}^{\mathsf{T}} R^{-1} \boldsymbol{z})}{(2\pi)^{M/2}(\det R)^{1/2}} \, \mathrm{d}\boldsymbol{z} &= \int_{\mathbb{R}^d} g(L(\boldsymbol{v}_{[1:M]})) \frac{\exp(-\frac{1}{2}\boldsymbol{v}^{\mathsf{T}} \boldsymbol{v})}{(2\pi)^{d/2}} \, \mathrm{d}\boldsymbol{v} \\
&= \int_{\mathbb{R}^d} g(L(\widetilde{L}^{-1}\boldsymbol{u})_{[1:M]}) \frac{\exp(-\frac{1}{2}\boldsymbol{u}^{\mathsf{T}} C^{-1} \boldsymbol{u})}{(2\pi)^{d/2}(\det C)^{1/2}} \, \mathrm{d}\boldsymbol{u} .
\end{aligned}$$

To complete the proof note that $L(\widetilde{L}^{-1}\boldsymbol{u})_{[1:M]} = (LL^{-1})\boldsymbol{u}_{[1:M]} = \boldsymbol{u}_{[1:M]}$. $\qquad\square$

**Lemma 2** *Let $C$ be as in Lemma 1 and suppose $C$ has a factorization of the form*

$$C = SS^{\mathsf{T}} . \tag{3.7}$$

*Then for any integrable function $f : \mathbb{R}^d \to \mathbb{R}$ we have*

$$\int_{\mathbb{R}^d} f(\boldsymbol{u}) \frac{\exp(-\frac{1}{2}\boldsymbol{u}^{\mathsf{T}} C^{-1} \boldsymbol{u})}{(2\pi)^{d/2}(\det C)^{1/2}} \, \mathrm{d}\boldsymbol{u} = \int_{[0,1]^d} f\big(S\boldsymbol{\Phi}_d^{-1}(\boldsymbol{x})\big) \, \mathrm{d}\boldsymbol{x} .$$

**Proof.** Straightforward, using the successive changes of variables $\boldsymbol{u} = S\boldsymbol{y}$ and $\boldsymbol{y} = \boldsymbol{\Phi}_d^{-1}(\boldsymbol{x})$. $\square$

It is easy to check that instead of the specific form (3.6), $C$ in Lemmas 1 and 2 may actually be any symmetric positive definite matrix chosen so that some $M \times M$ selection of $C$ coincides with $R$. By this we mean that selecting some subset of rows, and the same selection of columns, of $C$ gives us $R$. In this case, $\boldsymbol{u}_{[1:M]}$ in Lemma 1 is replaced by the vector $\boldsymbol{u}_R$ made up of the components of $\boldsymbol{u}$ corresponding to the particular selection of $C$.

**Corollary 3** *Let $C$ be any symmetric positive definite $d \times d$ matrix such that some $M \times M$ selection of $C$ coincides with $R$. If $C$ has the factorization (3.7), then (3.4) can be written as*

$$\mathbb{E}(\mathcal{G}_h(\boldsymbol{Z})) = \int_{[0,1]^d} \mathcal{G}_h\big((S\boldsymbol{\Phi}_d^{-1}(\boldsymbol{x}))_R\big) \, \mathrm{d}\boldsymbol{x} , \tag{3.8}$$

where, for any $\boldsymbol{u} \in \mathbb{R}^d$, $\boldsymbol{u}_R$ denotes the $M$ entries of $\boldsymbol{u}$ corresponding to the particular selection of $C$. The MC and QMC approximations for the expected value are, respectively,

$$\frac{1}{N} \sum_{n=1}^{N} \mathcal{G}_h\big((S\boldsymbol{y}^{(n)})_R\big) \qquad and \qquad \frac{1}{N} \sum_{n=1}^{N} \mathcal{G}_h\big((S\boldsymbol{\Phi}_d^{-1}(\boldsymbol{x}^{(n)}))_R\big), \tag{3.9}$$

where $\{\boldsymbol{y}^{(n)}\}$ are independent standard Gaussian random vectors from $\mathbb{R}^d$ and $\{\boldsymbol{x}^{(n)}\}$ are deterministically chosen QMC points from $[0,1]^d$.

## 3.3 Comparison with methods based on truncated KL expansion

A starting point for the solution of (2.1)–(2.4) is often taken to be the *Karhunen-Loève expansion*

$$Z(\vec{x}\,;\,\omega) \;=\; \sum_{\ell=1}^{\infty} \sqrt{\mu_\ell}\,\psi_\ell(\vec{x})\,Y_\ell(\omega), \quad \vec{x} \in D, \quad \omega \in \Omega\,, \tag{3.10}$$

where $Y_\ell(\omega)$ are independent standard Gaussian random variables, $\psi_\ell$ are orthonormal eigenfunctions (with respect to $L_2(D)$) and $\mu_\ell$ are the corresponding (positive) eigenvalues of the positive definite integral operator on $D$ with kernel function $r(\vec{x}, \vec{y})$. The sequence $\{\mu_\ell\}$ may be chosen non-increasing, and converges to 0 as $\ell \to \infty$.

For continuous covariance functions $r(\vec{x}, \vec{y})$ that are sufficiently well behaved near the diagonal $\vec{x} = \vec{y}$, the KL expansion (3.10) converges to a continuous function uniformly on $D$, with probability 1. In particular, Theorems 3.3.2 and 3.4.1 of [1] tell us that this convergence property holds for the covariance function (2.5). However, the convergence rate of (3.10) is determined by the decay of the eigenvalues $\mu_\ell$, and can be slow when $r$ is not smooth across the diagonal (as in (2.5)). In fact even in the 1-dimensional case the covariance function (2.5) has eigenvalues $\mu_\ell$ only of order $\ell^{-2}$. Therefore methods which attack (2.1)–(2.4) by first truncating (3.10) may suffer large truncation error unless the number of terms taken is large.

It is instructive to compare the approach described above with an approach based on the truncated KL expansion. Truncating (3.10) to $d$ terms and using it to evaluate the random vector $\boldsymbol{Z}$ given in (2.13) we obtain the approximation

$$\widetilde{\boldsymbol{Z}}(\omega) \;:=\; \Psi\boldsymbol{Y}(\omega)\,, \tag{3.11}$$

where $\boldsymbol{Y}(\omega)$ is a $d$-dimensional standard Gaussian random vector and $\Psi$ is the $M \times d$ matrix given by

$$\Psi_{i,\ell} \;=\; \sqrt{\mu_\ell}\,\psi_\ell(\vec{x}_i)\,, \quad i = 1, \ldots, M\,, \quad \ell = 1, \ldots, d\,. \tag{3.12}$$

From this, following the same arguments as above, we may approximate $\mathbb{E}(\mathcal{G}(Z))$ by

$$\mathbb{E}(\mathcal{G}(Z)) \stackrel{\mathrm{KL}}{\approx} \mathbb{E}(\mathcal{G}(\widetilde{\boldsymbol{Z}})) \stackrel{\mathrm{FE}}{\approx} \mathbb{E}(\mathcal{G}_h(\widetilde{\boldsymbol{Z}})) \;:=\; \int_{[0,1]^d} \mathcal{G}_h\big(\Psi\boldsymbol{\Phi}_d^{-1}(\boldsymbol{x})\big)\,\mathrm{d}\boldsymbol{x}. \tag{3.13}$$

MC and QMC approximations can then be obtained analogously to (3.9) replacing $S$ with $\Psi$.

The KL approach suffers from two drawbacks: the first is that the computation (approximation) of the KL eigenvalues and eigenfunctions can be very expensive, and the second is that there is an error in the truncated field (3.11) which is introduced before any finite element approximation is carried out.

## 3.4 Comparison with stochastic collocation

Suppose we have a general quadrature rule on the $d$-dimensional unit cube with weights $\{w_n\}$ and points $\{\boldsymbol{x}^{(n)}\}$:

$$\int_{[0,1]^d} f(\boldsymbol{x})\mathrm{d}\boldsymbol{x} \;\approx\; \sum_{n=1}^{N} w_n f(\boldsymbol{x}^{(n)}). \tag{3.14}$$

If this rule is applied to (3.13), then we obtain

$$\mathbb{E}(\mathcal{G}(Z)) \;\approx\; \sum_{n=1}^{N} w_n \mathcal{G}_h\big(\boldsymbol{\Psi}\boldsymbol{\Phi}_d^{-1}(\boldsymbol{x}^{(n)})\big)\,, \tag{3.15}$$

where $\boldsymbol{\Psi}$ is the matrix associated with the truncated KL expansion defined in (3.12). As above, the implementation of this requires the solution of $N$ different instances of (2.15). This is reminiscent of the stochastic collocation method of [2, 29, 30]. For example if (3.14) represents a quadrature rule based on applying sparse grid techniques to $d$-fold tensor products of one dimensional interpolatory rules, then the resulting approximation is essentially the method proposed in [29] (although the PDE to be solved in [29] is the primal instead of the mixed form (2.15)). An anisotropic variant of [29] which performs more work in the directions of the important variables and is more suitable for higher dimensions is given in [30]. However the quadrature rules in [29, 30] are different to the QMC rules which we study in this paper. In particular, our methods use equal weights and our quadrature points are more uniformly distributed in the $d$-dimensional unit cube. While the anisotropic sparse grid method and the QMC method both aim for robustness with respect to dimension, the dimensions of the examples which we report in this paper (up to $10^6$) go far beyond those tested in the examples in [30]. Moreover, as mentioned before, our method has the additional attraction that it does not involve truncation of the KL expansion.

The collocation methods can be seen as approximations of the older stochastic Galerkin method [3, 15, 33, 37], in which a variational approximation in both physical and probability space is performed, starting from a KL truncation of the random coefficient. These methods, together with stochastic collocation, have an extensive and elegant error analysis, e.g., [3, 37, 2, 29, 30]. The error analysis of our method will be one focus of our future work, but is not the aim of this paper.

## 4  Quasi-Monte Carlo methods

We begin this section with a brief but contemporary introduction to QMC theory, intended only to be sufficient for understanding the key issues concerning their use in §7.

### 4.1  QMC theory

Recall that QMC methods seek to approximate an integral of a function $f : [0,1]^d \to \mathbb{R}$ over the $d$-dimensional unit cube by an equal-weight quadrature rule

$$I_d(f) \;:=\; \int_{[0,1]^d} f(\boldsymbol{x})\,\mathrm{d}\boldsymbol{x} \quad \approx \quad \frac{1}{N}\sum_{n=1}^{N} f(\boldsymbol{x}^{(n)}) \;=:\; Q_{N,d}(f).$$
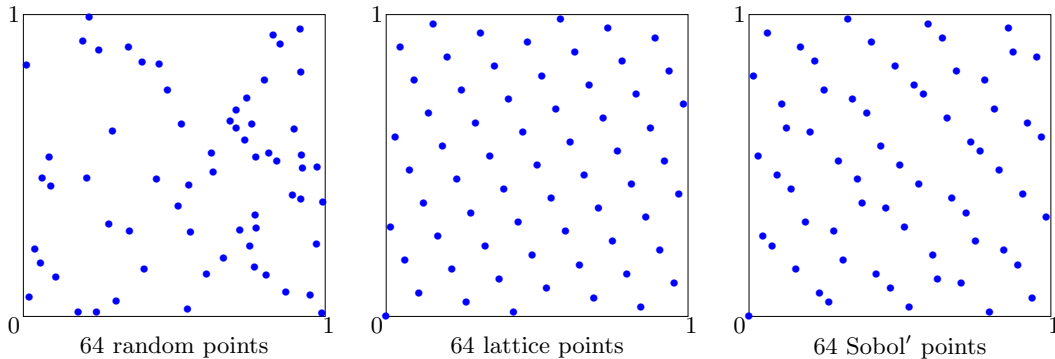
The name "quasi-Monte Carlo" reflects the fact that the simple Monte Carlo rule for the integral $I_d(f)$ has exactly the same appearance as $Q_{N,d}(f)$, albeit with a crucial difference, that in the MC rule the points $\{\boldsymbol{x}^{(n)}\}$ are chosen randomly and independently from a uniform distribution on $[0,1]^d$. The MC rule has the well known probabilistic error estimate

$$\sqrt{\mathbb{E}(I_d(f) - Q_{N,d}(f))^2} \;=\; \frac{\sigma(f)}{\sqrt{N}}, \qquad \text{where} \qquad \sigma^2(f) \;:=\; I_d(f^2) - (I_d(f))^2.$$

The general ambition of QMC rules is to improve upon the performance of the MC rule through a clever deterministic choice of points $\{\boldsymbol{x}^{(n)}\}$.

There are two main classes of QMC rules: so-called *lattice rules* and *nets*. Both classes of methods were initially introduced and studied by number theorists fifty years ago. Although we

Figure 1: Comparison of MC and QMC points



64 random points       64 lattice points       64 Sobol′ points

discuss examples of both classes of QMC rules here, our computations in §7 are restricted to *Sobol′ points* [41], which are the earliest and perhaps the most popular examples of nets.

In Figure 1 we show, for the 2-dimensional case, the first 64 Sobol′ points, compared to the points of a (good) 64-point lattice rule, and 64 (pseudo-) random points in the plane. Lattice points are very regular as seen in the figure: the lattice points form a group under addition modulo the integers, which includes all integer points if extended periodically from the unit square $[0,1]^2$ to all of $\mathbb{R}^2$. On the other hand, Sobol′ points are also very uniform albeit in a different sense: in every dyadic subdivision of the unit square into 64 identical pieces (for example, into 64 strips of size $1 \times \frac{1}{64}$, or 64 rectangles of size $\frac{1}{2} \times \frac{1}{32}$, etc), there is exactly 1 point. In general, Sobol′ point sets with $2^m$ points in $d$ dimensions are "$(t,m,d)$ nets" in base 2, meaning that each subdivision of the unit cube into $2^{m-t}$ identical pieces contains exactly $2^t$ points (as described above for $t = 0$, $m = 6$ and $d = 2$). We see that the uniform distribution of $(t,m,d)$ nets deteriorates as $t$ increases, and so we can think of $t$ as a "(lack of) quality parameter". Unfortunately, it is known that for Sobol′ points $t$ generally increases with $d$. This implies that Sobol′ points have an intrinsic order of importance of coordinate directions. They are more uniformly spread in the early dimensions and so the integration variables should always be ordered according to their importance.

The above description considers only Sobol′ point sets of size $2^m$ for some positive integer $m$, but practical Sobol′ point generators allow all intermediate vales of $N$ to be filled in, and indeed are usually presented as infinite sequences. For reviews of the classical material on QMC rules, see Niederreiter [27] and Sloan and Joe [38]. However, for both classes of QMC rules the justification for their use (and in the case of lattice rules, even the method of construction) has to be very different from the original arguments when $d$ is large: we think that either rule is better analyzed not in any of the classical function spaces, but rather in the *weighted* spaces introduced by Sloan and Woźniakowski [40].

The intuition behind the weighted spaces is that most problems are too hard when $d$ is large and all variables are of equal importance; but that the situation may be better if (with the variables appropriately ordered) the successive variables (that is, the successive components of $\boldsymbol{x}$ in $f(\boldsymbol{x})$) are of declining importance. To quantify the declining importance of successive variables, [40] introduced an infinite sequence of positive *weights* $\gamma_j$, satisfying

$$\gamma_1 \geq \gamma_2 \geq \cdots > 0.$$

These weights are then built into the function spaces in which the error analysis is done. Specifically, in its original version in [40] the $d$-dimensional weighted space $H_{d,\boldsymbol{\gamma}}$ (a real-valued Hilbert space) is a tensor product of 1-dimensional Hilbert spaces $H_{1,\gamma_j}$ with inner product

$$(f,g)_{1,\gamma_j} = f(1)\,g(1) + \frac{1}{\gamma_j} \int_0^1 f'(x)\,g'(x)\,\mathrm{d}x \tag{4.1}$$

11

and norm $\|f\|_{1,\gamma_j} = (f,f)_{1,\gamma_j}^{1/2}$. Other more general weighted spaces have been defined in the literature, but for our present purposes the details of weighted spaces are not important; only the weights themselves are important.

It is known from [18] that QMC rules exist for which

$$|I_d(f) - Q_{N,d}(f)| \leq C_\delta \, N^{-1+\delta} \, \|f\|_{H_{d,\gamma}}, \tag{4.2}$$

for every $\delta > 0$, and with $C_\delta$ independent of $d$ if

$$\sum_{j=1}^{\infty} \gamma_j^{1/2} < \infty. \tag{4.3}$$

A similar result under the weaker condition that the weights are summable was already given in [40], resulting in a convergence rate of $N^{-1/2}$. The importance of these results can be illustrated by noting that for the classical case, where all weights are equal, the constant grows exponentially with $d$.

In addition, we know how to construct "rank-1" lattice rules that achieve the error bound (4.2) under the condition (4.3) on the weights by means of the *component-by-component* (CBC) construction; see [22, 31, 39], or see [23] for a non-technical review of recent developments on weighted spaces and lattice rules. The points of a rank-1 lattice rule can be completely specified by one integer vector $\boldsymbol{g}$ whose components are all relatively prime to $N$. The essence of the CBC construction is that the components of $\boldsymbol{g}$ are determined one after the other, with the value of the $k$th component chosen such that it minimizes a certain quantity depending on the weights $\gamma_1, \ldots, \gamma_k$. Because the weights are decreasing and the earlier choices are in some sense more free than the later ones, the quality of the rule constructed by the CBC algorithm can be expected to deteriorate as the dimension increases — the same conclusion that we reached, for a different reason, for the Sobol′ sequence.

However, weighted spaces also provide a framework for understanding the success of the Sobol′ points in high dimensions. This comes from the work of Wang, [43], who showed that the Sobol′ points in any number of dimensions $d$ can achieve an error bound of the form (4.2) with $C_\delta$ independent of $d$, albeit with a condition on the weights that is slightly stronger than (4.3).

In summary, in the case of lattice rules the framework of weighted spaces provides an algorithm by which a suitable lattice rule can be constructed, once the weights are chosen. In the case of the Sobol′ points the weighted spaces (unlike the classical spaces) allow us to understand how and why the Sobol′ points (and indeed other nets) can be successful even in high dimensions. However, the key lesson we want the reader to take away from this section is that when using QMC rules it is crucial to order the components of the $d$-dimensional integration variable $\boldsymbol{x}$ in order of decreasing importance.

Parameters for obtaining good lattice or Sobol′ points can be found in `http://www.maths.unsw.edu.au/~fkuo`.

## 4.2   Randomization

It is well known that an estimate of the *standard error* for the MC method $Q_{N,d}(f)$ can be obtained from

$$s_N := \left( \frac{1}{N(N-1)} \sum_{n=1}^{N} \left( f(\boldsymbol{x}^{(n)}) - Q_{N,d}(f) \right)^2 \right)^{1/2}. \tag{4.4}$$

An empirical 95% *confidence interval* for the integral approximation (based on the assumption of Gaussian distribution of $Q_{N,d}(f) = \sum_{n=1}^{N} f(\boldsymbol{x}^{(n)})/N$) is then $Q_{N,d}(f) \pm 1.96 \, s_N$. Note that there is an equivalent formula for (4.4) which requires only one pass through the evaluations of $f$. But for both the one-pass and two-pass implementations, care must be taken to avoid numerical instability, see e.g., [19].

We recommend that QMC be implemented not in the pure, deterministic form described in §4.1, but rather in a randomized form that borrows from the MC method an estimate of the error from the spread of the results. There are two popular forms of randomization: *shifting* and *scrambling*, the former preserves the lattice structure while the latter preserves the net structure, see, e.g., [17]. Here we describe shifting and a simple form of scrambling called *digital shifting*.

We can add a *shift* $\boldsymbol{\Delta} \in [0,1]^d$ to any QMC point $\boldsymbol{x} \in [0,1]^d$, and "wrap" it back into the unit cube if necessary, to obtain a shifted point in the unit cube denoted by $\boldsymbol{x} \oplus \boldsymbol{\Delta} = \mathrm{frac}(\boldsymbol{x} + \boldsymbol{\Delta})$ where frac denotes the fractional part. This describes shifting. With digital shifting, we rewrite both $\boldsymbol{x}$ and $\boldsymbol{\Delta}$ in some base (2 is the natural choice for Sobol' points) and perform digitwise addition modulo the base. For example, with base 2, $\oplus$ is the bitwise exclusive-or operator: if $\boldsymbol{x} = (0.25, 0.75) = (0.01, 0.11)_2$ and $\boldsymbol{\Delta} = (0.375, 0.625) = (0.011, 0.101)_2$ then the digitally shifted point is $\boldsymbol{x} \oplus \boldsymbol{\Delta} = (0.001, 0.011)_2 = (0.125, 0.375)$.

To estimate the QMC error we make use of some number, say $\nu$, of random shifts (or digital shifts) of a selected QMC rule. Suppose the chosen QMC rule is a $\kappa$-point rule $Q_{\kappa,d}$ with points $\{\boldsymbol{x}^{(n)}\}$. Then the same rule shifted (or digitally shifted) by $\boldsymbol{\Delta} \in [0,1]^d$ is

$$Q_{\kappa,d,\boldsymbol{\Delta}}(f) := \frac{1}{\kappa} \sum_{n=1}^{\kappa} f(\boldsymbol{x}^{(n)} \oplus \boldsymbol{\Delta}),$$

and the randomly shifted (or digitally shifted) version of this rule with $N = \nu\kappa$ points is

$$Q_{N,d}^{\mathrm{ran}}(f) := \frac{1}{\nu} \sum_{i=1}^{\nu} Q_{\kappa,d,\boldsymbol{\Delta}^{(i)}}(f) = \frac{1}{\nu\kappa} \sum_{i=1}^{\nu} \sum_{n=1}^{\kappa} f(\boldsymbol{x}^{(n)} \oplus \boldsymbol{\Delta}^{(i)}), \qquad (4.5)$$

where $\boldsymbol{\Delta}^{(1)}, \ldots, \boldsymbol{\Delta}^{(\nu)}$ are $\nu$ independent random samples from a uniform distribution on $[0,1]^d$, such that the $Q_{\kappa,d,\boldsymbol{\Delta}^{(i)}}$ are independent estimates for $I_d(f)$. It is easy to verify that $Q_{N,d}^{\mathrm{ran}}(f)$ provides an *unbiased* estimate for $I_d(f)$. An approximation of the *standard error* for $Q_{N,d}^{\mathrm{ran}}(f)$ is given by

$$s_{N,\nu} := \left( \frac{1}{\nu(\nu-1)} \sum_{i=1}^{\nu} \left( Q_{\kappa,d,\boldsymbol{\Delta}^{(i)}}(f) - Q_{N,d}^{\mathrm{ran}}(f) \right)^2 \right)^{1/2}. \qquad (4.6)$$

To avoid spoiling the good QMC convergence rate, the number of random shifts (or digital shifts) should be small and fixed, say $\nu = 10$ or $20$. Then instead of taking the number of points $\kappa$ to be fixed, we increase $\kappa$ successively until the desired error threshold is satisfied. This is straightforward with the Sobol' sequence. A method to construct "extensible" (or embedded) lattice rules is described in [9].

# 5 Circulant embedding and FFT

To motivate this section, let us recall Corollary 3, which gives a formula for the expected value of $\mathcal{G}_h(\boldsymbol{Z})$, on the assumption that the covariance matrix $R$ of $\boldsymbol{Z}$ (recall (2.13)) can be embedded in a symmetric positive definite matrix $C$. In this section we show that $C$ can be taken to be a circulant matrix when the field is sampled at a regular rectangular array of points in $\mathbb{R}^2$. There are several new aspects to our discussion here compared to [5, 11]. In particular, the fact that the resulting $C$ can be diagonalized in real arithmetic is of importance to our algorithm.

To explain the circulant embedding technique we need to introduce some extra notation. Suppose that $A_0, \ldots, A_n$ are symmetric square matrices of the same size. Then the *symmetric block Toeplitz* matrix whose first block row and column contains these matrices is denoted by

$\mathrm{SBT}(A_0, \dots, A_n)$, i.e.,

$$\mathrm{SBT}(A_0, \dots, A_n) := \begin{bmatrix} A_0 & A_1 & A_2 & \cdots & A_n \\ A_1 & A_0 & A_1 & \cdots & A_{n-1} \\ A_2 & A_1 & A_0 & \cdots & A_{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ A_n & A_{n-1} & A_{n-2} & \cdots & A_0 \end{bmatrix}.$$

With this notation, the matrix $\mathrm{SBT}(A_0, \dots, A_n, A_{n-1}, \dots, A_1)$ is *block circulant*, that is, the $j$th block row of this matrix is just the same as the $(j-1)$th block row shifted one place to the right and wrapped around. Obviously these definitions also make sense in the special case when the blocks $A_i$ are scalars.

## 5.1 The 1D case

To facilitate understanding, it is useful to consider first the model case of a $1D$ random field $Z$, with covariance function $r(x, y) := \rho(|x - y|)$ for some $\rho : \mathbb{R}^+ \to \mathbb{R}$, and suppose this is to be sampled at $m+1$ equally spaced points (including the end points) in $[0, 1]$. (We remark that the absolute value in $\rho$ can be left out by some minor modifications.) Then the relevant covariance matrix is

$$R = \mathrm{SBT}(\rho_0, \dots, \rho_m), \quad \text{where} \quad \rho_j = \rho(j/m) \quad \text{for} \quad j = 0, \dots, m.$$

An embedding of $R$ in a $d \times d$ symmetric circulant matrix $C$ with $d := 2(m + J)$ can be obtained by taking

$$C = \mathrm{SBT}(\rho_0, \dots, \rho_m, \underbrace{\rho_{m+1}, \dots, \rho_{m+J}}_{\text{padding}}, \underbrace{\rho_{m+J-1}, \dots, \rho_1}_{\text{mirroring}}) \tag{5.1}$$

for some integer $J \geq 0$. Clearly $C$ contains $R$ in its top left-hand corner and the "mirroring" ensures the circulant structure for any choice of $J$. "Padding", i.e., taking $J \geq 1$ in (5.1), with large enough $J$ and suitable padding values, can ensure positive definiteness for all covariance functions of interest in this paper (see [6] and [11] for theory and numerical experiments). The choice $J = 0$, i.e., no padding, is sufficient for most of the examples in §7. When padding is needed, we specify our padding values by $\rho_j = \rho(j/m)$ for $j = m+1, \dots, m+J$. The matrix $C$ can then be written as

$$C = \mathrm{SBT}(\tilde{\rho}_0, \dots, \tilde{\rho}_{2(m+J)-1}), \quad \text{where} \quad \tilde{\rho}_j = \tilde{\rho}(j/m) \quad \text{for} \quad j = 0, \dots, 2(m+J) - 1,$$

and $\tilde{\rho} : \mathbb{R}^+ \to \mathbb{R}$ is the $2\ell$-periodic function (where $\ell := 1 + J/m$), defined by

$$\tilde{\rho}(t) := \begin{cases} \rho(t) & \text{for } t \in [0, \ell], \\ \rho(2\ell - t) & \text{for } t \in [\ell, 2\ell], \\ \tilde{\rho}(t - 2\ell) & \text{otherwise.} \end{cases}$$

The diagonalization of a $d \times d$ circulant matrix may be achieved in $\mathcal{O}(d \log d)$ operations by fast Fourier transform. Indeed, for a circulant matrix $C$ we have $C = F^{\mathsf{H}} \Lambda F$, where $F$ is the unitary Fourier matrix with entries

$$F_{p,q} = \frac{1}{\sqrt{d}} \exp(2\pi \mathrm{i}\, pq/d), \qquad p, q = 0, \dots, d - 1, \tag{5.2}$$

$F^{\mathsf{H}}$ is the Hermitian conjugate of $F$, and $\Lambda := \mathrm{diag}(\sqrt{d}\, F\, \boldsymbol{c})$ is the diagonal matrix of eigenvalues of $C$, with $\boldsymbol{c}$ denoting the first column of $C$. Note that since $C$ is symmetric and real and has real eigenvalues, we also have $C = F \Lambda F^{\mathsf{H}}$. The following lemma states that a real diagonalization of $C$ can be obtained by combining real and imaginary parts of $F$. This is sometimes called the (orthogonal) Hartley transform.

14

**Lemma 4** *Let $C$ be any symmetric circulant matrix and let $C = F\Lambda F^{\mathsf{H}}$. Then the matrix $G := \mathfrak{Re}(F) + \mathfrak{Im}(F)$ is real, symmetric and orthogonal, and $C = G\Lambda G^{\mathsf{T}}$.*

**Proof.** The matrices $F_c := \mathfrak{Re}(F)$ and $F_s := \mathfrak{Im}(F)$ are symmetric, and so is $G$. Since $F$ is unitary, we have $I = FF^H = (F_c + \mathrm{i}\,F_s)(F_c - \mathrm{i}\,F_s) = (F_c^2 + F_s^2) + \mathrm{i}\,(F_sF_c - F_cF_s)$. Hence

$$F_c^2 + F_s^2 = I \qquad \text{and} \qquad F_cF_s - F_sF_c = 0. \tag{5.3}$$

This and elementary trigonometry give

$$(F_sF_c)_{p,q} = (F_cF_s)_{p,q} = \frac{1}{d}\sum_{r=0}^{d-1}\cos\left(\frac{2\pi\,pr}{d}\right)\sin\left(\frac{2\pi\,rq}{d}\right) = 0.$$

Using (5.3) again we have $GG^{\mathsf{T}} = GG = F_c^2 + F_s^2 = I$, so $G$ is orthogonal.

Now, starting with $C = F\Lambda F^{\mathsf{H}}$, we multiply from the right by $F$ to obtain $CF = F\Lambda$. Separating this into real and imaginary parts yields $CF_c = F_c\Lambda$ and $CF_s = F_s\Lambda$. Summing the two equations and multiplying from the right by $G^{\mathsf{T}}$, we obtain the result. $\qquad\square$

## 5.2   The 2D case

Suppose $Z$ is a random field on a 2D domain with covariance function $r(\vec{x}, \vec{y}) = \rho(\|\vec{x} - \vec{y}\|_p)$, with $p = 1$ or $2$, and suppose we want to sample it on a uniform $(m_1 + 1) \times (m_2 + 1)$ rectilinear grid, with grid spacing $1/m_1$ in the $x_1$ direction and $1/m_2$ in the $x_2$ direction. Using a lexicographic ordering of points (first in the $x_1$ direction and then in the $x_2$ direction), the covariance matrix $R$ is *block Toeplitz with Toeplitz blocks* $R_0, \ldots, R_{m_2}$, i.e.,

$$R = \mathrm{SBT}(R_0, R_1, \ldots, R_{m_2}),$$

where each $R_j$ is itself an $(m_1 + 1) \times (m_1 + 1)$ Toeplitz matrix

$$R_j = \mathrm{SBT}(\rho_{0,j}, \rho_{1,j}, \ldots, \rho_{m_1,j}),$$

with entries $\rho_{i,j} := \rho(i/m_1, j/m_2)$, for $i = 0, \ldots, m_1$ and $j = 0, \ldots, m_2$.

We embed $R$ in a matrix $C$ which is *block circulant with circulant blocks*, by first embedding each Toeplitz block $R_j$ in a circulant matrix $C_j$ as in the 1D case, and then embedding the blocks in an analogous way. Again, padding may play a role in ensuring positive-definiteness. Thus, for padding parameters $J_1, J_2 \geq 0$, we set

$$C = \mathrm{SBT}(C_0, \ldots, C_{m_2}, \underbrace{C_{m_2+1} \ldots, C_{m_2+J_2}}_{\text{padding}}, \underbrace{C_{m_2+J_2-1}, \ldots, C_1}_{\text{mirroring}}),$$

where the blocks are themselves circulants

$$C_j = \mathrm{SBT}(\rho_{0,j}, \ldots, \rho_{m_1,j}, \underbrace{\rho_{m_1+1,j} \ldots, \rho_{m_1+J_1,j}}_{\text{padding}}, \underbrace{\rho_{m_1+J_1-1,j}, \ldots, \rho_{1,j}}_{\text{mirroring}}).$$

Analogous to the 1D case, we specify the padding by introducing the biperiodic function with period $2\ell_1$ in the $x_1$ direction and period $2\ell_2$ in the $x_2$ direction defined as

$$\tilde{\rho}(t_1, t_2) := \begin{cases} \rho(t_1, t_2) & \text{for } (t_1, t_2) \in [0, \ell_1] \times [0, \ell_2], \\ \rho(2\ell_1 - t_1, t_2) & \text{for } (t_1, t_2) \in [\ell_1, 2\ell_1] \times [0, \ell_2], \\ \tilde{\rho}(t_1, 2\ell_2 - t_2) & \text{for } (t_1, t_2) \in [0, 2\ell_1] \times [\ell_2, 2\ell_2], \\ \tilde{\rho}(t_1 - 2\ell_1, t_2 - 2\ell_2) & \text{otherwise,} \end{cases}$$

where $\ell_i = 1 + J_i/m_i$. The blocks of $C$ are chosen to be $C_j = \mathrm{SBT}(\tilde{\rho}_{0,j}, \ldots, \tilde{\rho}_{2(m_1+J_1)-1,j})$.

The matrix $C$ may be diagonalized by applying FFT techniques in each of the two coordinate directions in turn. This technique is well known and is described (in the context of a simpler matrix) in [42, p. 453–458] for example. To provide some more detail, let $d_i := 2(m_i + J_i)$, $i = 1, 2$. Then, $C$ is a $d \times d$ matrix with $d := d_1 d_2$, and we can identify the rows and columns of $C$ with the pairs $(p_1, p_2) \in \{0, \ldots, d_1 - 1\} \times \{0, \ldots, d_2 - 1\}$, arranged in lexicographical order. We then have $C = F \Lambda F^{\mathsf{H}}$, where $F$ is the $d \times d$ complex 2D Fourier matrix with entries

$$F_{(p_1,p_2),(q_1,q_2)} = \frac{1}{\sqrt{d}} \exp\left(\frac{2\pi \mathrm{i}\, p_1 q_1}{d_1}\right) \exp\left(\frac{2\pi \mathrm{i}\, p_2 q_2}{d_2}\right), \quad p_i, q_i = 0, \ldots, d_i - 1, \quad i = 1, 2.$$

Here, as in 1D, $\boldsymbol{c}$ denotes the first column of $C$, and $\Lambda = \sqrt{d}\,\mathrm{diag}(F\boldsymbol{c})$ contains the eigenvalues of $C$. The following result is the 2D analogue of Lemma 4.

**Lemma 5** *Let $d = d_1 \times d_2$ and let $C$ be any $d \times d$ symmetric block circulant matrix of the form $C = \mathrm{SBT}(C_0, \ldots, C_{d_2-1})$ with $d_1 \times d_1$ symmetric circulant blocks $C_j$, $j = 0, \ldots, d_2 - 1$. If $C = F \Lambda F^{\mathsf{H}}$ then $G := \mathfrak{Re}(F) + \mathfrak{Im}(F)$ is a real, symmetric and orthogonal matrix, and $C = G \Lambda G^{\mathsf{T}}$.*

**Proof.** It is easy to verify for the 2D case that $FF^{\mathsf{H}} = I$ and that $\mathfrak{Re}(F)\,\mathfrak{Im}(F) = 0$. The proof then follows the same lines as the proof of Lemma 4. □

## 5.3 Randomized QMC approximation combined with circulant embedding

We now outline the steps required for implementing a randomized version (see §4.2) of the QMC approximation in (3.9) combined with circulant embedding when the vector $\boldsymbol{Z} \in \mathbb{R}^M$ represents the random field $Z$ sampled at a regular rectangular array of points. From Lemma 5 we see that in the factorisation (3.7) of $C$ the matrix $S$ can be chosen to be $S = G\Lambda^{1/2}$.

As we explained in §4.1, the opening dimensions of QMC points are of higher quality, and for this reason we need to identify the importance of the integration variables and relabel them accordingly. It turns out that the magnitude of the eigenvalues of the circulant $C$ provides a good guide, which is why we permute the QMC components in accordance with the non-increasing sorted eigenvalues (see Step (iii) of the preprocessing below).

**Preprocessing:** (i) Embed $R$ into a $d \times d$ matrix $C$ that is block circulant with circulant blocks, and thus obtain $\boldsymbol{c}$, the first column of $C$. (ii) Compute $\boldsymbol{\lambda} := \sqrt{d}\,\mathtt{fft}(\boldsymbol{c})$ using a 2D FFT routine, e.g., FFTW [14]. (iii) Check that all eigenvalues in $\boldsymbol{\lambda}$ are real and positive. If not, increase the padding in $C$ (thus increasing $d$ by 1) and go back to (i). (iv) Choose a permutation $\pi$ such that $(\lambda_{\pi(j)})_{j=1}^d$ are in non-increasing order.

**For each QMC point $\boldsymbol{x}^{(n)}$ and each shift $\boldsymbol{\Delta}^{(i)}$:** (i) Compute $\boldsymbol{y} := \boldsymbol{\Phi}_d^{-1}(\boldsymbol{x}^{(n)} \oplus \boldsymbol{\Delta}^{(i)})$ where $\oplus$ denotes either shifting or digital shifting. (ii) Evaluate $\boldsymbol{w} := (\sqrt{\lambda_j}\, y_{\pi^{-1}(j)})_{j=1}^d$. (iii) Compute $\boldsymbol{v} := \mathtt{fft}(\boldsymbol{w})$. (iv) Take $\boldsymbol{u} := \mathfrak{Re}(\boldsymbol{v}) + \mathfrak{Im}(\boldsymbol{v})$. (v) Set $\boldsymbol{Z} := (\boldsymbol{u})_R$ as described in Corollary 3. (vi) Solve (2.15) with this instance of $\boldsymbol{Z}$.

Thus, with $\kappa$ denoting the number of QMC points and $\nu$ denoting the number of shifts, we obtain $N = \nu\kappa$ realizations of $\boldsymbol{Z}$ and $N$ solutions of (2.15). The expected value and its error estimate are computed as in (4.5) and (4.6).

## 5.4 Comparison to the MC approach

We finish this section by highlighting some of the differences between our QMC adaptation and the MC strategy described, for example, in [11]. Recall that QMC methods take the point of view of integration, not that of simulation, and therefore the justification for the embedding
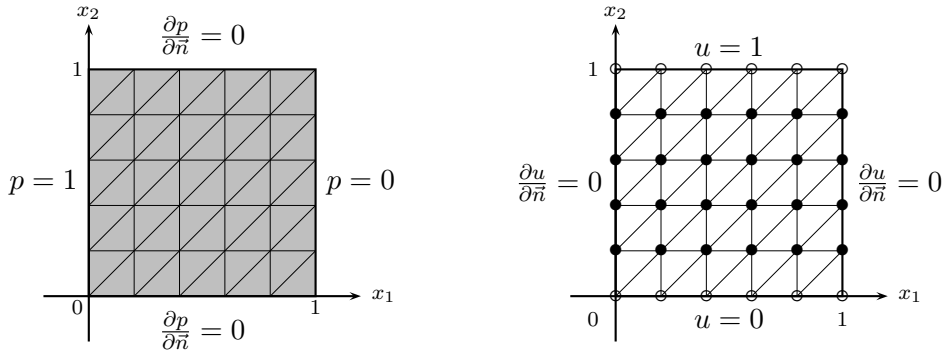
Figure 2: Left: Model problem and typical grid. Right: Boundary conditions and grid for auxiliary problem (6.9)–(6.10).

approach is quite different, see §3.2. While the MC error can be estimated from the spread of all function evaluations (since each sample is independent from the others), to estimate the QMC error we need randomization, see §4.2, which has the added advantage of removing bias from the QMC estimate. Whereas QMC points need to be mapped from the unit cube back to $\mathbb{R}^d$ using the inverse of the cumulative normal distribution function $\boldsymbol{\Phi}_d^{-1}$, the MC method can avoid this transformation by generating normal variates directly, and thus saving some computational work. However, as shown in §7, this is only a minor part of the total cost.

The permutation of variables as described in §5.3 is important for QMC methods, but it has no effect on the performance of the MC method because the MC samples only need to have the correct mean and covariance. For the same reason, the MC method has extra freedom in extracting $M$ components from $S\boldsymbol{y}^{(n)}$ (see (3.9)) in the sense that a different selection of components could be used for each simulation.

Finally we note that the MC strategy in [11] makes use of the complex factorization of $C$ to transform *one* complex vector of length $d$, i.e., *two* real input vectors, to obtain *two* real output vectors by separating the real and imaginary parts of the transformed vector. However, for QMC methods it is crucial that the assignment of QMC components to the integration variables remains the same for every function evaluation, which is why we use the real factorization in terms of $G$ in Lemmas 4 and 5. Application of the real factorization has the same computational cost as that of the complex factorization using standard FFT tricks.

## 6  Fast implementation of mixed finite elements

In this section we restrict to the case when (2.1)–(2.2) are to be solved on the simple domain $D = (0,1)^2$ subject to the specific mixed conditions

$$p(0, x_2) = 1 \quad \text{and} \quad p(1, x_2) = 0 \quad \text{for all} \quad x_2 \in [0,1], \quad \text{and} \tag{6.1}$$

$$q_2(x_1, 0) = 0 \quad \text{and} \quad q_2(x_1, 1) = 0 \quad \text{for all} \quad x_1 \in [0,1]. \tag{6.2}$$

This is sometimes referred to as a *flow cell* in the literature (cf. [25]). We show how to efficiently compute three quantities of physical interest using Raviart-Thomas mixed finite elements (FEs) on the uniform triangular meshes with $m$ subdivisions in each coordinate direction depicted in Figure 2 (left). For convenience we denote $h = 1/m$. We use the divergence-free reduction technique introduced in [7, 35] (see also [12]) to solve the resulting saddle point systems (2.11) or (2.15). The good computation times in §7 depend crucially on this fast procedure.

## 6.1 Three important physical quantities

The simplest quantity of interest that we will study is the **pressure head** $p(\vec{x}^*; \cdot)$ at a given point $\vec{x}^* \in D$.

The second quantity of interest is the **effective permeability** of a block of random porous media (cf. [25]), which in the case of the simple flow cell $D = (0, 1)^2$ with the boundary conditions (6.1)–(6.2) is given by

$$k_{\text{eff}}(\omega) := \frac{\int_D q_1(\vec{x}; \omega) \, d\vec{x}}{\int_D -\frac{\partial p}{\partial x_1}(\vec{x}; \omega) \, d\vec{x}} \, .$$

It measures the mean behavior of the random porous medium in $D$ and is of interest, e.g., in stochastic homogenization. Note that on the simple flow cell $D = (0, 1)^2$ we have $\int_D -\frac{\partial p}{\partial x_1} \, d\vec{x} = \int_0^1 (p(0, x_2) - p(1, x_2)) \, dx_2 = 1$.

Finally, in the context of modeling underground waste repositories, it is important to study the motion of pollutant particles in the velocity field $\vec{q}$, because transport by flowing groundwater is the main mechanism for pollutants (such as radionuclides) to return to man's environment (see e.g. [28]). For each realization of the permeability field $k$, we place a particle at, say $\vec{x}^\dagger$, and follow its path through the field $\vec{q}$ until it exits the domain $D$. The final position and the time it takes to reach the boundary, the so-called **breakthrough time**, are important quantities of physical interest and in §7, we shall compute expected values of the breakthrough time $T = T(\omega)$ for the flow cell $D = (0, 1)^2$ and $\vec{x}^\dagger = (0, 0.5)^\mathsf{T}$.

## 6.2 Divergence-free reduction of mixed finite element systems

It can be shown that (2.15) has a unique solution, but it is highly ill conditioned when $h$ is small and when $k$ varies strongly. The fact that it is indefinite further complicates the direct application of fast multilevel solvers, such as algebraic multigrid (AMG). This is the reason why we first use an algebraic reduction of (2.15) to a symmetric positive definite system and a triangular system. To explain this reduction we define

$$\mathring{\mathcal{V}}_h := \{\vec{v} \in \mathcal{V}_h : b(w, \vec{v}) = 0 \quad \text{for all } w \in \mathcal{W}_h\} \, ,$$

the subspace of (discretely) divergence-free Raviart-Thomas elements. Since (2.15) has a unique solution, we have $\dim \mathring{\mathcal{V}}_h = n_v - n_w =: \mathring{n}$ and we can consider choosing the first $\mathring{n}$ elements $\vec{v}_1, \ldots, \vec{v}_{\mathring{n}}$ in the basis of $\mathcal{V}_h$, such that they also form a basis of $\mathring{\mathcal{V}}_h$. Then $B_{i,\ell} = b(w_\ell, \vec{v}_i) = 0$ for $i = 1, \ldots, \mathring{n}$, and since the solution $\vec{q}_h$ belongs to $\mathring{\mathcal{V}}_h$, we have $Q_{\mathring{n}+1} = \cdots = Q_{n_v} = 0$. Thus, we can decouple the first block-row of (2.15) into two independent (square) problems for $\mathring{\mathbf{Q}} := (Q_1, \ldots, Q_{\mathring{n}})^\mathsf{T}$ and $\mathbf{P}$, namely

$$\sum_{j=1}^{\mathring{n}} \widetilde{M}_{i,j} Q_j = g_i \qquad \qquad \text{for all } i = 1, \ldots, \mathring{n}, \qquad (6.3)$$

$$\sum_{\ell=1}^{n_w} B_{\mathring{n}+s,\ell} P_\ell = g_{\mathring{n}+s} - \sum_{j=1}^{\mathring{n}} \widetilde{M}_{\mathring{n}+s,j} Q_j \qquad \text{for all } s = 1, \ldots, n_w. \qquad (6.4)$$

We will see below that $\vec{v}_{\mathring{n}+1}, \ldots, \vec{v}_{n_v}$ can be chosen such that the matrix $(B_{\mathring{n}+s,\ell})_{1 \leq \ell, s \leq n_w}$ is bi-diagonal and (6.4) can therefore be solved in $\mathcal{O}(n_w)$ operations by simple back substitution. The core task is to solve (6.3). This system is 5 times smaller than the original system (2.15) (cf. [7]) and simpler to solve as we will explain now.

It is a property of the de Rham complex, which describes the connection between various FE spaces including the Raviart-Thomas elements (cf. [4, 35]), that a basis for $\mathring{\mathcal{V}}_h$ can be constructed from curls of the standard, scalar-valued, continuous, piecewise linear FE basis associated with $\mathcal{T}_h$. To be more precise, for any node $\vec{x}_j$ of $\mathcal{T}_h$ let $\varphi_j$ be the piecewise linear hat function associated with $\vec{x}_j$ such that $\varphi_j(\vec{x}_{j'}) = \delta_{j,j'}$, for all nodes $\vec{x}_{j'}$ of $\mathcal{T}_h$. Further, let us partition the set $\mathcal{N}$ of all nodes $\vec{x}_j$ of $\mathcal{T}_h$, into a set $\mathcal{N}_B$ of nodes on the bottom boundary where $x_2 = 0$, a set $\mathcal{N}_T$ of nodes on the top boundary where $x_2 = 1$, and the remainder which we call $\mathcal{N}_I$. Since degrees of

freedom in $\mathcal{V}_h$ can be associated with edges of $\mathcal{T}_h$ and those in $\mathcal{W}_h$ with elements, it follows from Euler's polyhedron theorem that $\mathring{n} = \#\mathcal{N}_I + 1$ (cf. [7, Theorem 4.3]). Let $\vec{x}_1, \ldots, \vec{x}_{\mathring{n}-1} \in \mathcal{N}_I$. Then we can choose the basis for $\mathring{\mathcal{V}}_h$ in the following way (cf. [7, §4.2]):

$$
\begin{aligned}
\vec{v}_i &:= \vec{\mathrm{curl}}\,\varphi_i\,, \quad i = 1, \ldots, \mathring{n} - 1, \quad \text{and} \\
\vec{v}_{\mathring{n}} &:= \vec{\mathrm{curl}}\,\varphi_T, \quad \text{where} \quad \varphi_T := \sum_{\vec{x}_j \in \mathcal{N}_T} \varphi_j,
\end{aligned}
\tag{6.5}
$$

and where $\vec{\mathrm{curl}}f := (\partial f/\partial x_2, -\partial f/\partial x_1)^{\mathsf{T}}$ in 2D. (Note that the choice of the top, rather than the bottom boundary in (6.5) is not essential.) Substituting (6.5) into (2.14) we have

$$
\widetilde{M}_{i,j} = \int_D \bar{k}^{-1}\,\vec{\mathrm{curl}}\,\varphi_j \cdot \vec{\mathrm{curl}}\,\varphi_i\,\mathrm{d}\vec{x} = \underbrace{\int_D \bar{k}^{-1}\,\vec{\nabla}\varphi_j \cdot \vec{\nabla}\varphi_i\,\mathrm{d}\vec{x}}_{=:\,a(\varphi_i, \varphi_j)}, \quad i, j = 1, \ldots, \mathring{n} - 1, \tag{6.6}
$$

where $\bar{k}$ is a piecewise constant function defined elementwise in (2.14). Also from (2.9) we have

$$
g_i = G(\vec{v}_i) = -\int_{\Gamma_{\mathrm{in}}} \vec{\mathrm{curl}}\,\varphi_i \cdot \vec{n}\,\mathrm{d}\Gamma(\vec{x}) = 0, \quad i = 1, \ldots, \mathring{n} - 1,
$$

with $\Gamma_{\mathrm{in}}$ denoting the left-hand boundary $\{0\} \times [0, 1]$. Similarly, $\widetilde{M}_{i,\mathring{n}} = a(\varphi_i, \varphi_T)$, $\widetilde{M}_{\mathring{n},\mathring{n}} = a(\varphi_T, \varphi_T)$ and $g_{\mathring{n}} = 1$. Now denote

$$
A_{i,j} := a(\varphi_i, \varphi_j) \quad \text{and} \quad f_i := -a(\varphi_i, \varphi_T), \quad \text{for} \ \ 1 \le i, j \le \mathring{n} - 1.
$$

Then the solution of (6.3) can be found via block–factorization (see [7] for details) to be

$$
\mathring{\mathbf{Q}} = \eta_h \begin{pmatrix} A^{-1}\mathbf{f} \\ 1 \end{pmatrix}, \quad \text{where} \quad \eta_h := \frac{1}{\widetilde{M}_{\mathring{n},\mathring{n}} - \mathbf{f}^{\mathsf{T}}A^{-1}\mathbf{f}}. \tag{6.7}
$$

Hence, solving (6.3) reduces to solving

$$
A\,\mathbf{u} = \mathbf{f} \quad \text{and then using} \quad \mathring{\mathbf{Q}} = \eta_h \begin{pmatrix} \mathbf{u} \\ 1 \end{pmatrix}. \tag{6.8}
$$

Now $A\mathbf{u} = \mathbf{f}$ is a standard (continuous, piecewise linear) FE discretization of the auxiliary PDE

$$
\vec{\nabla} \cdot (k^{-1}\,\vec{\nabla}u) = 0 \tag{6.9}
$$

subject to the mixed boundary conditions

$$
\begin{aligned}
u(x_1, 0) = 0 \quad &\text{and} \quad u(x_1, 1) = 1 \quad \text{for all} \ \ x_1 \in [0, 1], \quad \text{and} \\
\tfrac{\partial u}{\partial x_1}(0, x_2) = 0 \quad &\text{and} \quad \tfrac{\partial u}{\partial x_1}(1, x_2) = 0 \quad \text{for all} \ \ x_2 \in [0, 1],
\end{aligned}
\tag{6.10}
$$

i.e., a second-order elliptic problem of the same form as that obtained for the pressure $p$ by substituting (2.1) into (2.2), but with diffusion coefficient $k^{-1}$ instead of $k$ and the roles of Dirichlet and Neumann boundary interchanged (see Figure 2, right, for an illustration). The auxiliary function $u$ plays the role of a stream function for the divergence-free vector field $\vec{q}$ in the original system of PDEs (2.1), (2.2). We emphasise again however, that $\bar{k}$ in the definition of the stiffness matrix $A$ in (6.6) depends only on the values of the random field contained in the vector $\mathbf{Z} \in \mathbb{R}^M$.

There are a number of robust solvers for (6.8) and some recent theory analyzing the resilience of many of the methods to large variations in the diffusion coefficient (see, e.g., [16, 32]). In this paper we will use algebraic multigrid (AMG) to solve (6.8), and we will see that its cost grows linearly with the problem size and is not affected by variations in $k$. AMG has been shown experimentally to be a very efficient method for solving high contrast diffusion problems but unfortunately there is no theoretical justification of its coefficient robustness yet.
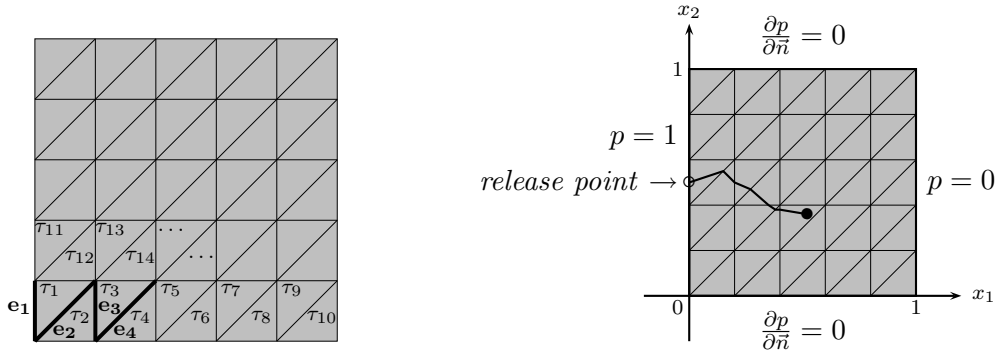
Figure 3: Left: Element numbering and distinguished edge $e_\ell$ associated with element $\tau_\ell$. Right: Piecewise linear particle path (particle released at $(0, 0.5)^\mathsf{T}$)

## 6.3 Simple derivation of quantities of interest

We will show now that all the quantities of interest defined above can be computed in a straightforward way from the solution $\mathbf{u}$ of the auxiliary problem (6.8). Note first that using (6.5), (6.7) and (6.8), and setting $u_i = 1$ for all $\vec{x}_i \in \mathcal{N}_T$ and $u_i = 0$ for all $\vec{x}_i \in \mathcal{N}_B$ (cf. (6.10)) we have

$$\vec{q}_h \;=\; \sum_{i=1}^{\mathring{n}} Q_i \vec{v}_i \;=\; \eta_h \left( \sum_{\vec{x}_i \in \mathcal{N}_I} u_i \,\vec{\mathrm{curl}}\,\varphi_i \;+\; \sum_{\vec{x}_i \in \mathcal{N}_T} \vec{\mathrm{curl}}\,\varphi_i \right) \;=\; \eta_h \sum_{\vec{x}_i \in \mathcal{N}} u_i \,\vec{\mathrm{curl}}\,\varphi_i \,. \qquad (6.11)$$

Now, since the functions $\varphi_i$ are piecewise linear w.r.t. $\mathcal{T}_h$, $\vec{q}_h$ is piecewise constant. We denote the values of $\mathbf{u}$ locally, at the nodes of each triangle $\tau \in \mathcal{T}_h$, by $u_{\mathrm{NW}(\tau)}$, $u_{\mathrm{NE}(\tau)}$, $u_{\mathrm{SW}(\tau)}$, if $\tau$ is above the diagonal, and $u_{\mathrm{NE}(\tau)}$, $u_{\mathrm{SW}(\tau)}$ $u_{\mathrm{SE}(\tau)}$, if $\tau$ is below the diagonal, where NW, NE, SW and SE stand for north-west, north-east, south-west and south-east, respectively. Recall that $h = 1/m$. Then it is easily seen that

$$\vec{q}_h|_\tau \;=\; \begin{cases} \dfrac{\eta_h}{h} \left( u_{\mathrm{NW}(\tau)} - u_{\mathrm{SW}(\tau)} \,,\; u_{\mathrm{NW}(\tau)} - u_{\mathrm{NE}(\tau)} \right)^\mathsf{T} & \text{if } \tau \text{ is above the diagonal,} \\[2ex] \dfrac{\eta_h}{h} \left( u_{\mathrm{NE}(\tau)} - u_{\mathrm{SE}(\tau)} \,,\; u_{\mathrm{SW}(\tau)} - u_{\mathrm{SE}(\tau)} \right)^\mathsf{T} & \text{if } \tau \text{ is below the diagonal.} \end{cases} \qquad (6.12)$$

It turns out that implicitly, in the divergence-free reduction, we have already computed the FE approximation $k_{\mathrm{eff},h}(\omega) := \int_D q_{h,1}(\vec{x}\,;\,\omega)\,\mathrm{d}\vec{x}$ of the effective permeability $k_{\mathrm{eff}}$.

**Proposition 6** Let $\eta_h$ be as defined in (6.7). Then $k_{\mathrm{eff},h} = \eta_h$.

**Proof.** For any $i \in \{0, \ldots, m\}$ define $\bar{q}_{h,1}(i) := \int_0^1 q_{h,1}(ih, x_2)\,\mathrm{d}x_2$. Then, since $\vec{q}_h$ is (discrete) divergence-free in the Raviart–Thomas case, i.e., $\int_\tau \vec{\nabla}\cdot\vec{q}_h\,\mathrm{d}\vec{x} = 0$ for all $\tau \in \mathcal{T}_h$, and since $q_{h,2} = 0$ at the top and bottom boundary (cf. (6.2)), it follows from the Divergence Theorem, applied elementwise on the rectangle $(ih, jh) \times (0, 1)$, that

$$\bar{q}_{h,1}(j) \;-\; \bar{q}_{h,1}(i) \;=\; \int_{(ih, jh) \times (0,1)} \nabla \cdot \vec{q}_h\,\mathrm{d}\vec{x} \;=\; 0\,, \quad \text{for all} \quad 0 \le i \le j \le m\,. \qquad (6.13)$$

Thus $\bar{q}_{h,1}(i)$ is constant (as a function of $i$) and so $\int_D q_{h,1}\,\mathrm{d}\vec{x} = \int_0^1 q_{h,1}(0, x_2)\,\mathrm{d}x_2$. Hence, using the first part of (6.12)

$$k_{\mathrm{eff},h} \;=\; h \sum_{\tau:\,\tau \cap \Gamma_{\mathrm{in}} \neq \emptyset} \frac{\eta_h}{h} \left( u_{\mathrm{NW}(\tau)} - u_{\mathrm{SW}(\tau)} \right) \;=\; \eta_h(1 - 0)\,.$$

$\square$

To recover the coefficients of the piecewise constant approximation of the pressure $p_h$ via (6.4), we need a complementary basis $\{\vec{v}_{\mathring{n}+1}, \ldots, \vec{v}_{n_v}\} \subset \mathcal{V}_h$ to $\{\vec{v}_1, \ldots, \vec{v}_{\mathring{n}}\}$ in (6.5). Since $n_v - \mathring{n} = n_w$, we need exactly one basis function per element $\tau \in \mathcal{T}_h$. It turns out that these can be chosen from amongst the standard set of basis functions for Raviart–Thomas elements associated with edges of $\mathcal{T}_h$. Let us order the elements $\tau_1, \ldots, \tau_{n_w} \in \mathcal{T}_h$ as depicted in Figure 3 (left), i.e., starting with the bottom row, numbering the elements from left to right and then proceeding with the second row, etc. For each triangle $\tau_\ell$ we also choose a distinguished edge $e_\ell$ as shown, i.e., the vertical edge for $\ell$ odd and the diagonal edge for $\ell$ even. The standard Raviart-Thomas basis function $\vec{v}_{\mathring{n}+\ell} \in \mathcal{V}_h$ associated with edge $e_\ell$ is uniquely defined by

$$\int_{e'} \vec{v}_{\mathring{n}+\ell} \cdot \vec{n}_{e'} \, \mathrm{d}s \;=\; \delta_{e_\ell, e'}, \qquad \text{for all edges } e' \text{ of } \mathcal{T}_h, \tag{6.14}$$

where $\vec{n}_{e'}$ is the unit normal on edge $e'$ chosen to lie in $\{\vec{x} : x_1 > 0\} \cup \{(0,1)^\mathsf{T}\}$. It can be shown that with this set of complementary basis functions, the functions $\vec{v}_1, \ldots, \vec{v}_{n_v}$ are linearly independent and thus form a basis for $\mathcal{V}_h$ (cf. [35, 36] for details). Thus, with the usual basis $\{w_\ell\}$ of $\mathcal{W}_h$ given by the characteristic functions of the elements $\tau_\ell$ of $\mathcal{T}_h$, it follows from the divergence theorem and (6.14) that

$$B_{\mathring{n}+s,\ell} = b(w_\ell, \vec{v}_{\mathring{n}+s}) = -\int_{\tau_\ell} \vec{\nabla} \cdot \vec{v}_{\mathring{n}+s} \, \mathrm{d}\vec{x} = \begin{cases} 1, & \text{if } \ell = s, \\ -1, & \text{if } \ell = s+1 \text{ and } \tau_\ell \cap \mathrm{supp}(\vec{v}_{\mathring{n}+s}) \neq \emptyset, \\ 0 & \text{otherwise.} \end{cases}$$

Similarly

$$g_{\mathring{n}+s} \;=\; G(\vec{v}_{\mathring{n}+s}) \;=\; -\int_{\Gamma_{\mathrm{in}}} \vec{v}_{\mathring{n}+s} \cdot \vec{n} \, \mathrm{d}\Gamma(\vec{x}) \;=\; \begin{cases} 1, & \text{if } \tau_s \cap \Gamma_{\mathrm{in}} \neq \emptyset, \\ 0 & \text{otherwise.} \end{cases}$$

Substituting these last two equations into (6.4) we get the simple recursion formula

$$P_\ell \;=\; \begin{cases} 1 & - \Delta P_\ell, & \text{if } \tau_\ell \cap \Gamma_{\mathrm{in}} \neq \emptyset, \\ P_{\ell-1} & - \Delta P_\ell, & \text{otherwise,} \end{cases} \tag{6.15}$$

where $\Delta P_\ell := (\widetilde{M}\mathbf{Q})_{\mathring{n}+\ell}$. It turns out that each $\Delta P_\ell$ can be computed in a simple way from (6.12) on the two elements adjacent to edge $e_\ell$. We omit the details. Since $p_h$ is piecewise constant with respect to the triangulation $\mathcal{T}_h$, its value is not defined on any of the vertices or edges of the triangulation. In §7 we will study the expected value of the pressure at the center of the flow cell, and so we simply extend the definition of $p_h$ to all of $D$ by averaging:

$$p_h(\vec{x}) \;:=\; \begin{cases} p_h(\vec{x}) & \text{if } \vec{x} \in \mathrm{int}(\tau), \ \tau \in \mathcal{T}_h & \text{(interior of elements),} \\ \frac{1}{6} \sum_{\tau:\vec{x}\in\tau} p_h|_\tau & \text{if } \vec{x} \in \mathcal{N}_I \backslash \partial D & \text{(interior vertices),} \\ \frac{1}{2} \sum_{\tau:\vec{x}\in\tau} p_h|_\tau & \text{for all other } \vec{x} \in D & \text{(interior edges).} \end{cases} \tag{6.16}$$

The following proposition shows an exactness result for the mean pressure which turns out to be useful in §7 for checking the performance of our code.

**Proposition 7** *Let $D = (0,1)^2$ and consider the problem (2.1)–(2.2) subject to (6.1)–(6.2), with finite element approximation as described in §2.2. Then*

$$\mathbb{E}(p_h(\vec{x}^*\,;\,\cdot)) \;=\; \mathbb{E}(p(\vec{x}^*\,;\,\cdot)) \;=\; 1/2, \qquad \text{when} \qquad \vec{x}^* = (1/2, 1/2)^\mathsf{T}.$$

**Proof.** Introduce the bijection $\gamma(\vec{x}) := (1-x_1, 1-x_2)^\mathsf{T}$ on $D$, and define $k_\gamma(\vec{x}\,;\,\omega) := k(\gamma(\vec{x})\,;\,\omega)$, $\vec{q}_\gamma(\vec{x}\,;\,\omega) := \vec{q}(\gamma(\vec{x})\,;\,\omega)$ and $p_\gamma(\vec{x}\,;\,\omega) := 1 - p(\gamma(\vec{x})\,;\,\omega)$. Then $(\vec{q}_\gamma, p_\gamma)$ satisfy (2.1)–(2.2) with random field $k_\gamma$ instead of $k$ but with the same boundary conditions (6.1)–(6.2). Since $k = \log Z$

21

and $Z$ is a mean-zero Gaussian random field with a covariance function $r(\vec{x}, \vec{y})$ that is invariant under the transformation $\vec{x} \to \gamma(\vec{x})$, we can deduce that $k_\gamma$ is equal in law to $k$ and thus

$$\mathbb{E}(p(\vec{x}\,;\,\cdot)) \;=\; \mathbb{E}(p_\gamma(\vec{x}\,;\,\cdot)) \;=\; \mathbb{E}(1 - p(\gamma(\vec{x})\,;\,\cdot)) \;=\; 1 - \mathbb{E}(p(\gamma(\vec{x})\,;\,\cdot))\,,$$

and putting $\vec{x} = \vec{x}^*$ yields $\mathbb{E}(p(\vec{x}^*\,;\,\cdot)) = 1/2$.

To complete the proof note that the grid $\mathcal{T}_h$ is invariant under the transformation $\vec{x} \to \gamma(\vec{x})$. Hence, we can as in the continuous case define FE functions $\vec{q}_{h,\gamma}(\vec{x}\,;\,\omega) := \vec{q}_h(\gamma(\vec{x})\,;\,\omega)$ and $p_{h,\gamma}(\vec{x}\,;\,\omega) := 1 - p_h(\gamma(\vec{x})\,;\,\omega)$ that solve the mixed FE system (2.10) with $k_\gamma$ instead of $k$, and deduce in a similar way that $\mathbb{E}(p_h(\vec{x}^*\,;\,\cdot)) = 1/2$. $\qquad\square$

We make use of Proposition 7 in §7, when we study the error in expected value of pressure at $(1/2, 1/2)^\mathsf{T}$, since in this special case there is no discretization error and so we can study the convergence of the QMC methods in isolation.

Since $\vec{q}_h$ is piecewise constant w.r.t. $\mathcal{T}_h$, particle paths and travel times can trivially be computed in an element by element fashion (cf. Figure 3, right). In particular, if a particle enters a triangle $\tau \in \mathcal{T}_h$ through edge $e^{\mathrm{in}}$ at point $\vec{x}^{\mathrm{in}}$, the travel time to any of the other two edges $e' \subset \tau$ is

$$\frac{\mathrm{dist}(\vec{x}^{\mathrm{in}}, e')}{\vec{q}_h \cdot \vec{n}_\tau|_{e'}}.$$

Since $\vec{q}_h$ has continuous normal components across all edges $e$ of $\mathcal{T}_h$, the travel time to one of the two edges $e' \neq e^{\mathrm{in}}$ has to be positive. If the travel time to both edges is positive, then the actual travel time in $\tau$ is the minimum of the two. If they are equal, the particle exits $\tau$ through a vertex and we arbitrarily choose one of the edges.

Again the travel time of a particle in $\tau$ and the position of the point where it leaves $\tau$ depend only on the terms on the right hand side of (6.12), as well as on some geometric considerations. By following the particle path from element to element through the domain $D$ and summing up the travel times in each element, it follows that the finite element approximation $T_h(\omega)$ of the breakthrough time $T(\omega)$ is just some nonlinear function of the differences of the values of $\mathbf{u}$ at neighboring grid points.

# 7  Numerical results

In this section we examine in detail the performance of the MC and randomized QMC methods for computing the expected value of each of the three physical quantities defined in §6, namely, pressure head at the point $\vec{x}^* = (1/2, 1/2)^\mathsf{T}$, effective permeability, and breakthrough time. The algorithms for computing these quantities using the mixed finite element (FE) method were described in §6, but for convenience we sample the random field $Z$ at the midpoints of the diagonal edges of the grid and use a different quadrature rule in the assembly of $\widetilde{M}$ in (2.14), namely the hypothenuse rule such that here $\overline{k}_\tau := \exp(Z(\vec{m}_\tau\,;\,\omega))$ in (2.14) with $\vec{m}_\tau$ being the midpoint of the diagonal edge of $\tau$. We use the FFTW package [14] for the circulant embedding technique and the AMG1R5 code [34] to solve the large, ill-conditioned linear equation systems (6.8) for each individual realization. We take digitally shifted Sobol$'$ points as our representative for randomized QMC methods. To generate Sobol$'$ points we use the parameters from [21].

Our numerical results are all in 2D, for the covariance function (2.5) with "1-norm" ($p = 1$) and "2-norm" ($p = 2$). We study a range of parameters $\sigma$ and $\lambda$ in (2.5). These are given in Table 1, where generally the cases on the right are more difficult than those on the left.

In each case we are computing expected values of our quantities of interest of the general form (3.8), approximated using mixed FEs on a uniform grid with $m$ subdivisions in each coordinate direction. Again we set $h = 1/m$. The resulting high-dimensional integral is transformed to the domain $[0, 1]^d$. The transformed integrand and its dimensionality $d$ both depend on $m$, typically

| Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|---|---|---|---|---|
| $\sigma^2 = 1,\ \lambda = 1$ | $\sigma^2 = 1,\ \lambda = 0.3$ | $\sigma^2 = 1,\ \lambda = 0.1$ | $\sigma^2 = 3,\ \lambda = 1$ | $\sigma^2 = 3,\ \lambda = 0.1$ |

Table 1: Five choices of $\lambda$ and $\sigma$ which we consider (difficulty goes from left to right)

$d = \mathcal{O}(m^2)$. The largest dimension considered is $d \approx 2 \times 10^6$ . To be more precise, referring to the notation in §5.2, we have taken $m_1 = m_2 = m$, subdivided the domain $D$ into $m^2$ equal squares and then divided these into triangles by drawing the diagonals from bottom left to top right. The field $Z$ is sampled at the $m^2$ centres of these squares. The dimension of the matrix $R$ is $M = m^2$, and the dimension of the matrix $C$ is $d = 4(m-1)^2$ when no padding is used.

There are in general two main contributions to the error: the discretization error, caused by the FE discretization in space, and the quadrature error of the MC or the randomized QMC method. To estimate the MC quadrature error we use the standard error estimate $s_N$, see (4.4), with increasing number of samples $N$. To estimate the randomized QMC quadrature error we use the standard error estimate $s_{N,\nu}$, see (4.6), with $\nu = 16$ random digital shifts for Sobol' points. Here $N$ is the total number of sample points, i.e., the number of QMC points multiplied by 16. To estimate the discretization error we apply our method on a sequence of grids of decreasing mesh width $h$, using a sufficiently large value of $N$, and deduce empirically an error estimate from these results via linear regression.

## 7.1 Timing tests

We start by looking at the computational cost. There is some initial cost to factorize the circulant matrix $C$ but this only requires one FFT application and is negligible (less than 1% of the total time for more than 30 samples, even on our finest mesh, see Table 2). For each individual sample the computational time is made up of one FFT application, one application of the (vectorized) inverse $\mathbf{\Phi}_d^{-1}$ of the cumulative standard normal distribution function, and one linear solve with AMG1R5. In Table 2 we give timings for the individual parts for 1000 samples for Case 1. All the CPU-times quoted in this section were obtained on a standard laptop with a 2GHz Intel T7300 processor.

| $m$ | $d$ | Setup | Applying $\mathbf{\Phi}_d^{-1}$ | FFT | AMG | Total |
|---|---|---|---|---|---|---|
| 33 | 4096 | 0.00 | 1.0 (17%) | 0.22 (4%) | 4.5 (76%) | 5.9 |
| 65 | 16384 | 0.01 | 3.9 (17%) | 1.2 (5%) | 16.5 (75%) | 22 |
| 129 | 65536 | 0.06 | 15 (16%) | 5.1 (6%) | 67 (73%) | 92 |
| 257 | 262144 | 0.15 | 62 (16%) | 31 (8%) | 290 (73%) | 400 |
| 513 | 1048576 | 0.6 | 258 (15%) | 145 (8%) | 1280 (73%) | 1750 |
|  | $\mathcal{O}(m^2)$ | $\mathcal{O}(m^2)$ | $\mathcal{O}(m^2)$ | $\mathcal{O}(m^2 \log m^2)$ | $\sim \mathcal{O}(m^2)$ | $\sim \mathcal{O}(m^2)$ |

Table 2: Timings (in seconds) for 1000 samples in our randomized QMC implementation on a standard laptop (2GHz Intel T7300) for Case 1 with 1-norm covariance function. (Percents of the total time in brackets.)

We see that the majority of the cost comes from the linear solves ($\sim 75\%$). The cost of the linear solve and the application of $\mathbf{\Phi}_d^{-1}$ grow like $\mathcal{O}(m^2)$ as the mesh is refined, which is optimal. The FFT part grows like $\mathcal{O}(m^2 \log m^2)$, but even on the finest mesh which we consider its cost is the smallest of the three ($\sim 8\%$). The cost for our MC implementation is identical to that of randomized QMC. It could be slightly reduced if the normal variates are generated directly to avoid the application of $\mathbf{\Phi}_d^{-1}$, but the cost of the linear solve, which is the dominant part, would always remain the same.

Let us study the effect of changing $\lambda$ and $\sigma^2$. The cost of applying $\mathbf{\Phi}_d^{-1}$ as well as the FFT are not affected. The cost for AMG grows slightly with $\sigma^2 \to \infty$ and with $\lambda \to 0$, but even in the hardest case that we study here (i.e., Case 5) it only grows by 40–50%. For example, the corresponding times in Table 2 for Case 5 with $m = 257$ are $t_{\text{setup}} = 0.13$, $t_{\mathbf{\Phi}_d^{-1}} = 64$, $t_{\text{FFT}} = 31$, $t_{\text{AMG}} = 440$ and $t_{\text{total}} = 550$.

Changing from the 1-norm to the 2-norm covariance function does have a substantial effect for correlation lengths $\lambda > 0.2$, since it is necessary in that case to introduce padding to ensure the circulant embedding is positive definite (see §5; padding is not necessary for the 1-norm covariance). As a consequence the dimension $d$ may be substantially larger in the 2-norm case. For example, with $\lambda = 1$ and $m = 65$ we need $d \approx 5.6 \times 10^5$, and the corresponding times in Table 2 are $t_{\mathbf{\Phi}_d^{-1}} = 132$ and $t_{\text{FFT}} = 75$, while $t_{\text{AMG}} = 16$ as in the 1-norm case. For $\lambda \le 0.2$ no padding is necessary even in the 2-norm case and the timings are similar to those in Table 2.

## 7.2 Pressure head at centre

We now analyze the convergence for the expected value of the pressure head at $\vec{x}^* = (1/2, 1/2)^{\mathsf{T}}$. This is an ideal test case to study the convergence of the quadrature error as $N$ increases, since $\mathbb{E}(p_h(\vec{x}^*)) = \mathbb{E}(p(\vec{x}^*)) = 0.5$ for any $h$ and the discretization error is 0 (see Proposition 7).

In Figure 4 we plot the standard error estimates $s_N$ and $s_{N,16}$ (for MC and randomized QMC, respectively) when approximating $\mathbb{E}(p_h(\vec{x}^*))$ in each of the Cases 1–5. In all graphs the QMC results are marked by (blue) stars while the (green) crosses are the MC results. The two triangles illustrate $\mathcal{O}(N^{-1})$ and $\mathcal{O}(N^{-1/2})$ convergence. We see clearly the superior performance of the QMC method over the MC method. The convergence rates above each graph are computed by linear regression. As expected, MC converges with a rate of $\mathcal{O}(N^{-1/2})$ throughout. The convergence rate of QMC is clearly better than that of MC, even though it does degenerate somewhat in the harder cases (i.e., when $\lambda$ is smaller or $\sigma^2$ larger). As we see from the bottom four plots in Figure 4, the mesh size $m$ and hence the dimension $d$ of the quadrature domain does not seem to have any influence on the convergence rates.

In absolute terms the increased convergence rate of QMC means that in Case 4, for example, for $m = 129$ to achieve an error of $10^{-3}$ we only require $N = 2500$ samples in QMC while $N = 45000$ samples are necessary in MC. In CPU-time this difference equates to 4 minutes versus 1.25 hours. The difference is even more dramatic if we seek an error of $10^{-4}$: for $m = 129$, while $N = 35000$ samples suffice in QMC, more than 4.5 million samples are necessary in MC, which equates to 1 hour for QMC versus more than 5 days for MC.

## 7.3 Effective permeability

Analogous plots for the standard errors in the expected value of the effective permeability $\mathbb{E}(k_{\text{eff},h})$ are given in Figure 5. Again the QMC convergence rate decreases slightly as the problem becomes harder, but, as in the case of the pressure head, QMC requires many fewer samples than MC to achieve a desired accuracy. For example taking again Case 4 with $m = 129$, it suffices to choose $N = 250000$ to achieve an accuracy of $10^{-3}$ in QMC, while MC would require almost 30 million samples, which is 7 hours versus about 1 month. Again the convergence rates do not seem to be influenced by the dimension $d$.

However, in contrast to the pressure head computations, the expected value of $k_{\text{eff},h}$ varies as $h$ changes due to discretization error. To quantify the discretization error and to decide, given a required tolerance $\varepsilon$, which mesh size to use, we estimate $\mathbb{E}(k_{\text{eff},h})$ as accurately as possible using QMC with $N = 2.1$ million samples on a sequence of grids with $h \to 0$. Estimates for $\mathbb{E}(k_{\text{eff},h})$, together with 95% confidence intervals (obtained by $\pm 1.96\, s_{N,16}$), in Cases 1, 2, 3 and 5 are presented in Table 3.

If we make an assumption that the discretization error for $\mathbb{E}(k_{\text{eff},h})$ decays like $\mathcal{O}(h^\beta)$, then the value of $\beta$ and the exact value $\mathbb{E}(k_{\text{eff}})$ can be estimated numerically via linear regression.
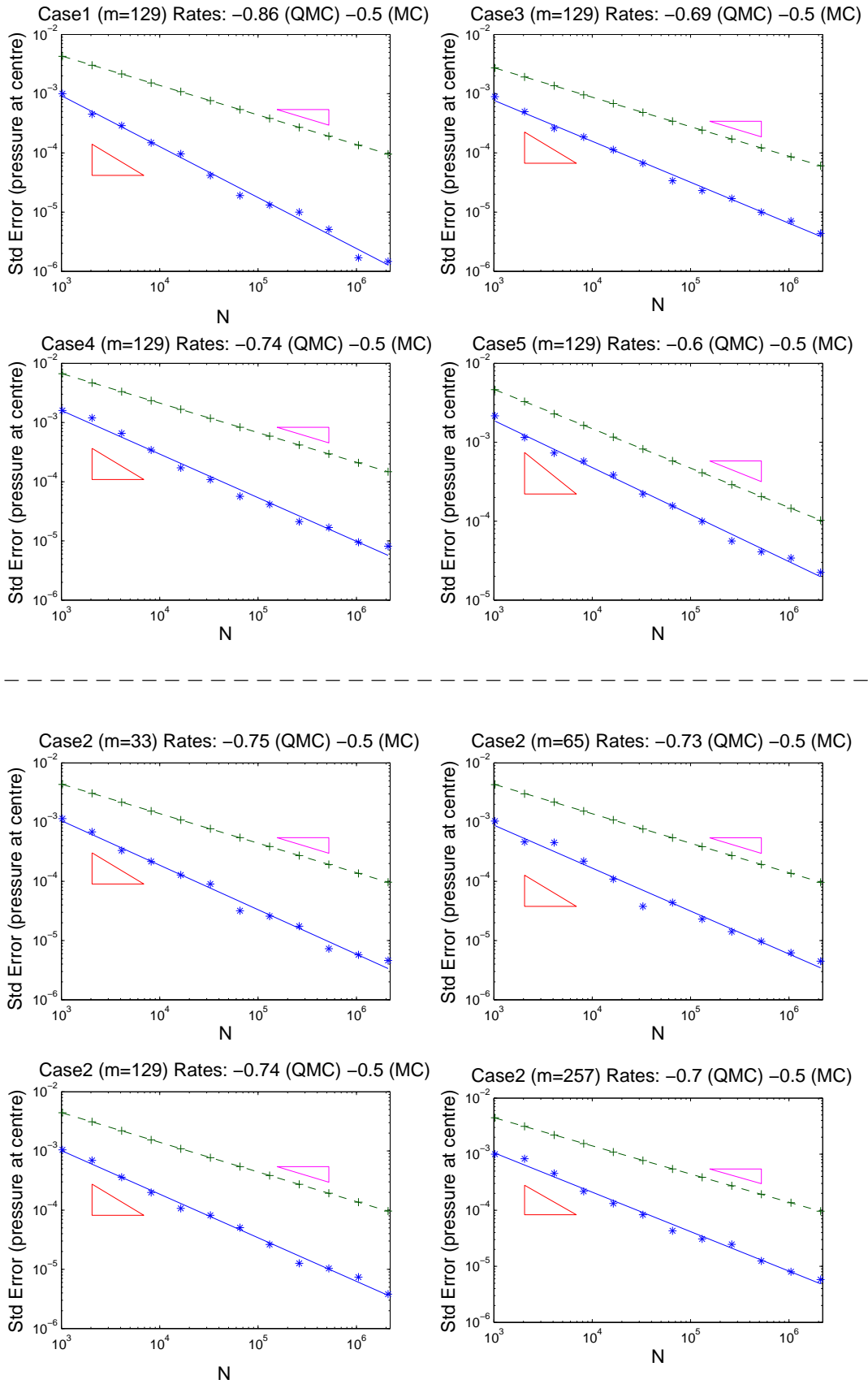
Figure 4: Standard errors $s_N$ and $s_{N,16}$ of expected pressure at centre for the 1-norm covariance function (i.e., $p = 1$ in (2.5)): Cases 1, 3, 4, 5 (above, with $m = 129$) and Case 2 (below, for various values of $m$).
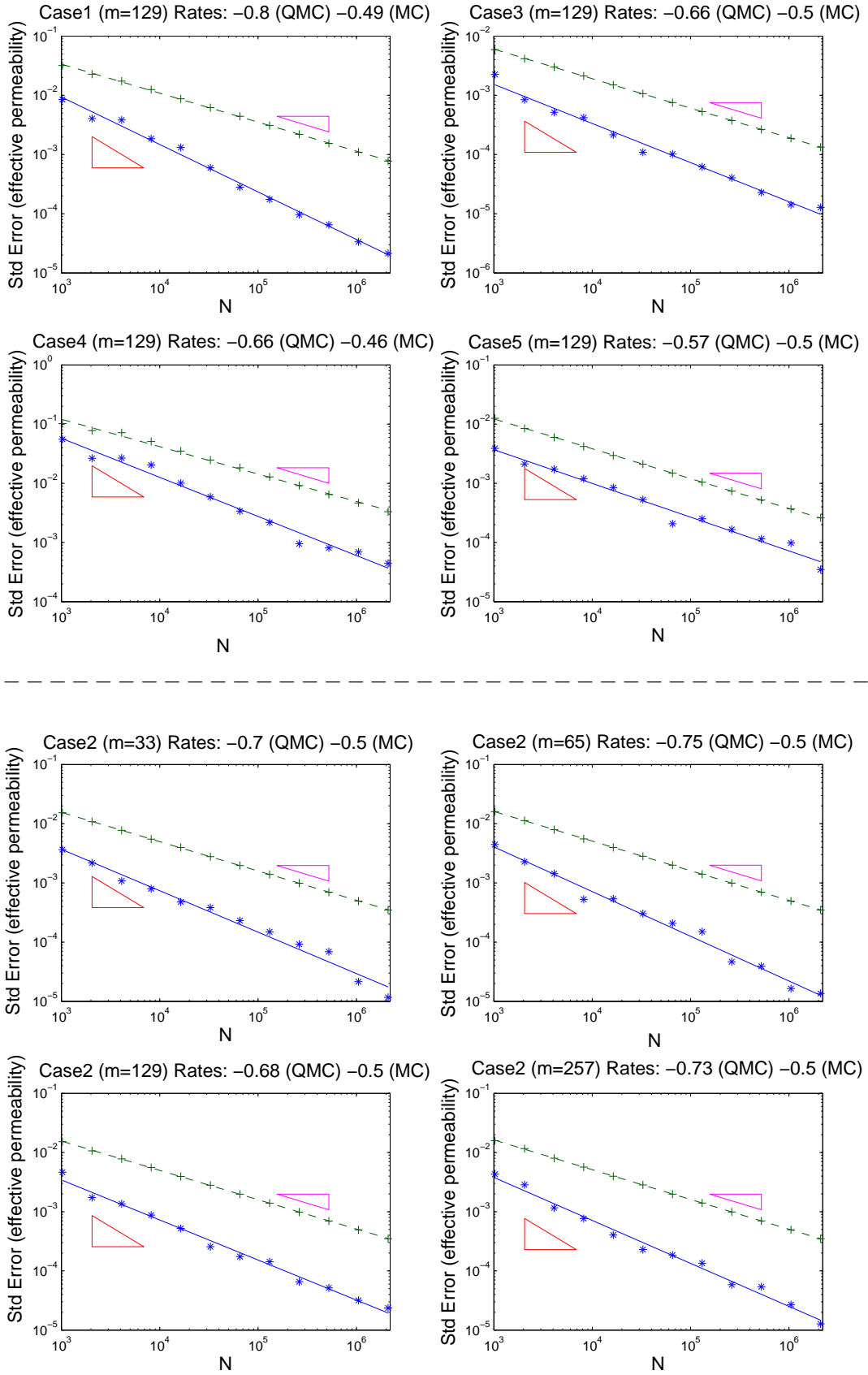
Figure 5: Standard errors $s_N$ and $s_{N,16}$ of expected effective permeability $k_{\mathrm{eff},h}$ for the 1-norm covariance function (i.e., $p = 1$ in (2.5)): Cases 1, 3, 4, 5 (above, for $m = 129$) and Case 2 (below, for various values of $m$).

| $1/h$ | Case 1 | Case 2 | Case 3 | Case 5 |
|---|---|---|---|---|
| 33 | $1.314970 \pm 3.5[\text{-5}]$ | $1.097452 \pm 2.3[\text{-5}]$ | $1.000958 \pm 2.7[\text{-5}]$ | $1.022453 \pm 8.6[\text{-5}]$ |
| 65 | $1.315102 \pm 4.1[\text{-5}]$ | $1.099467 \pm 2.7[\text{-5}]$ | $1.011073 \pm 2.0[\text{-5}]$ | $1.045101 \pm 1.1[\text{-4}]$ |
| 129 | $1.315174 \pm 4.2[\text{-5}]$ | $1.100222 \pm 4.6[\text{-5}]$ | $1.015403 \pm 2.5[\text{-5}]$ | $1.055474 \pm 6.8[\text{-5}]$ |
| 257 | $1.315234 \pm 3.6[\text{-5}]$ | $1.100459 \pm 2.5[\text{-5}]$ | $1.017014 \pm 2.4[\text{-5}]$ | $1.059588 \pm 1.4[\text{-4}]$ |
| $\infty$ | $1.315215 \pm 1.2[\text{-5}]$ | $1.100811 \pm 1.1[\text{-4}]$ | $1.018492 \pm 1.9[\text{-4}]$ | $1.062633 \pm 3.1[\text{-4}]$ |

Table 3: Estimates for $\mathbb{E}(k_{\mathrm{eff},h})$ computed with QMC (with $N \approx 2.1 \times 10^6$). The confidence intervals are estimated using 16 random shifts. In the last row we give estimates for the exact value of $\mathbb{E}(k_{\mathrm{eff}})$ obtained via linear regression (discarding $h = 1/257$ in Case 1 since the standard error is bigger than the FE error for $N \approx 2.1 \times 10^6$ in this case).

Numerical experimentation shows that the optimal $\beta$ is about 1.25 in each of the cases. The estimate for $\mathbb{E}(k_{\mathrm{eff}})$ using $\beta = 1.25$ (with its 95% confidence interval) is given in the last row of Table 3 for each case. In Table 4 we list our estimates for the discretization error for Cases 1, 2, 3 and 5. Note that the behavior of the error is indeed very close to $\mathcal{O}(h^{1.25})$.

| $1/h$ | Case 1 | Case 2 | Case 3 | Case 5 |
|---|---|---|---|---|
| 33 | $0.000245 \pm 3.7[\text{-5}]$ | $0.003358 \pm 1.1[\text{-4}]$ | $0.017533 \pm 2.0[\text{-4}]$ | $0.040179 \pm 3.2[\text{-4}]$ |
| 65 | $0.000113 \pm 4.3[\text{-5}]$ | $0.001344 \pm 1.1[\text{-4}]$ | $0.007419 \pm 1.9[\text{-4}]$ | $0.017531 \pm 3.3[\text{-4}]$ |
| 129 | $0.000041 \pm 4.4[\text{-5}]$ | $0.000589 \pm 1.2[\text{-4}]$ | $0.003089 \pm 2.0[\text{-4}]$ | $0.007158 \pm 3.2[\text{-4}]$ |
| 257 | $0.000019 \pm 3.8[\text{-5}]$ | $0.000352 \pm 1.1[\text{-4}]$ | $0.001477 \pm 1.9[\text{-4}]$ | $0.003044 \pm 3.4[\text{-4}]$ |

Table 4: Estimates for $\mathbb{E}(k_{\mathrm{eff}} - k_{\mathrm{eff},h})$ computed by linear regression on the results in Table 3 using $\beta = 1.25$ (discarding $h = 1/257$ in Case 1).

We can also see from Table 4 that in each of the cases we can expect an accuracy of at most $\mathcal{O}(10^{-4})$ for $\mathbb{E}(k_{\mathrm{eff},h})$ on any computationally feasible mesh. Recall that 1000 samples on a grid with mesh size $h = 1/513$ require about 30–40 minutes on a single processor (see Table 2). The mesh size grows very rapidly with $\lambda \to 0$, and so for Cases 3 and 5 where $\lambda = 0.1$, only an accuracy of about $\mathcal{O}(10^{-3})$ is possible.

In Table 5 we list the mesh sizes necessary to attain a discretization error smaller than $10^{-3}$ in Cases 1, 2, 3 and 5, together with the number of samples necessary for QMC or MC to ensure that the 95% confidence level of $\mathbb{E}(k_{\mathrm{eff},h})$ (i.e., $1.96\,s_{N,16}$ for QMC and $1.96\,s_N$ for MC) is also smaller than $10^{-3}$. Interestingly the number of samples necessary to achieve a fixed tolerance decreases with $\lambda \to 0$ (Cases 1–3), but this is offset by the larger discretization error for smaller $\lambda$, so that Cases 2 and 3 are indeed harder than Case 1 when measured in CPU time (see Table 5). Case 5 is as expected the hardest. We see that $\sigma^2$ not only influences the discretization error (see Table 4), it also influences the QMC and the MC error and so a substantially larger number of samples is required in Case 5 than in Case 4. Again we see clearly the advantage of QMC over MC, with a computational time of 28 minutes versus 28 hours in Case 2, and similar speedups for the other cases.

## 7.4 Breakthrough time

For the breakthrough time $T(\omega)$, the advantage of QMC over MC is less pronounced as we can see in Figure 6, but it still outperforms MC consistently with convergence rates of about $N^{-0.57}$, and a significantly smaller constant. The discretization error is similar to that of $\mathbb{E}(k_{\mathrm{eff},h})$ for the cases where $\sigma^2 = 1$, but for $\sigma^2 = 3$ it is substantially larger. Concentrating only on Cases 2 and 5 and proceeding as above we obtain the results in Table 6 (with $\beta = 1.35$).
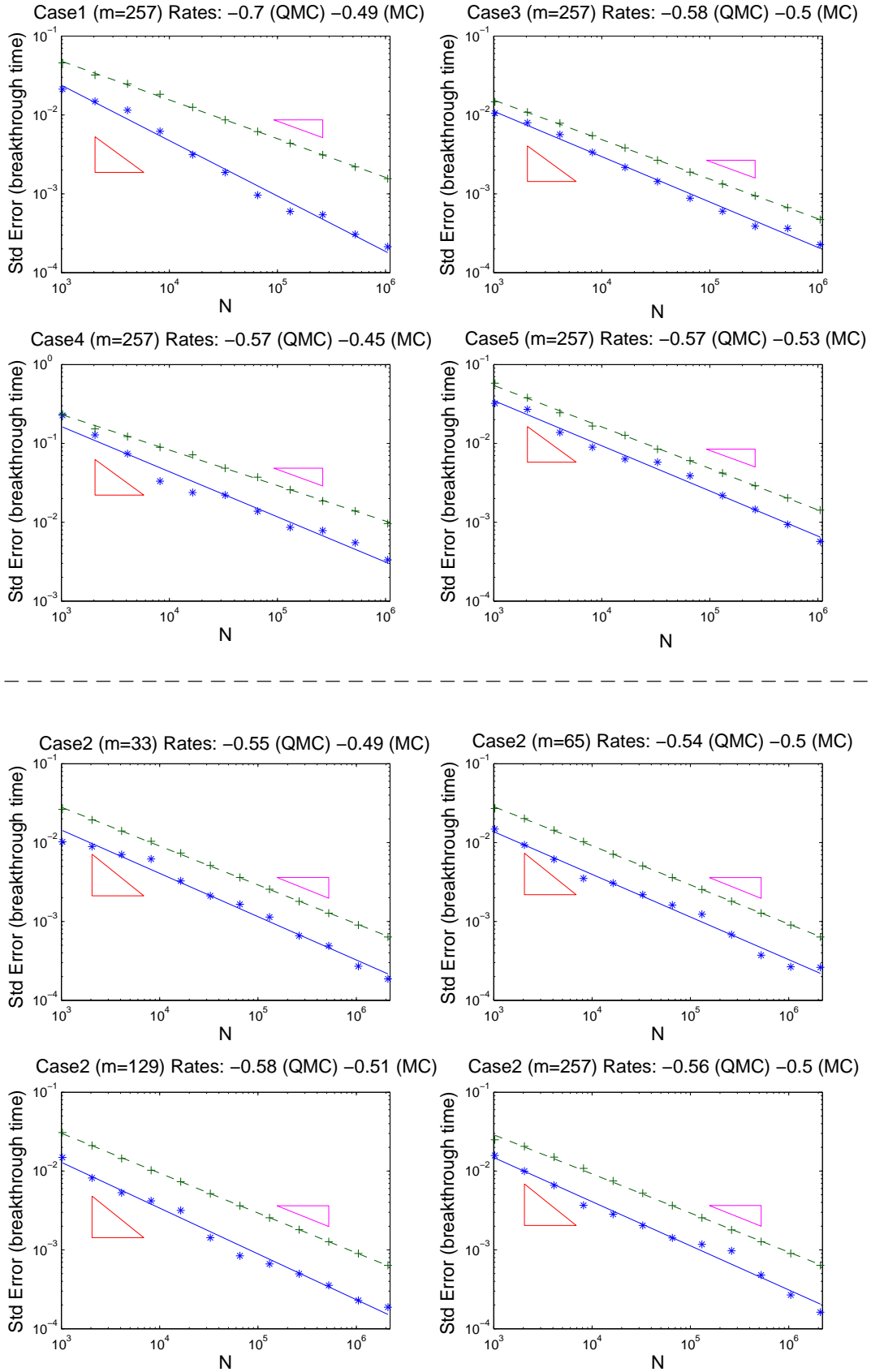
Figure 6: Standard errors $s_N$ and $s_{N,16}$ of expected breakthrough time $T_h$ for the 1-norm covariance function (i.e., $p = 1$ in (2.5)): Cases 1, 3, 4, 5 (above, $m = 257$) and Case 2 (below, for various values of $m$).

| Case | $1/h$ | # Samples $N$ | | CPU Time | | | | $\mathbb{E}(k_{\text{eff},h})$ |
|------|-------|------|---------|-----------|------|----|------|-------------------|
| | | QMC | MC | QMC | | MC | | |
| 1 | 17 | 64000 | 4250000 | 1.5 min | | 1.75 h | | 1.3147 |
| 2 | 129 | 16700 | 982000 | 28 min | | 28 h | | 1.1003 |
| 3 | 513 | 4900 | 142000 | 160 min | | 80 h | $\approx$ 3 d | 1.0182 |
| 5 | 1025 | 28700 | 525000 | 2100 min | $\approx$ 35 h | 1600 h | $\approx$ 67 d | 1.0614 |

Table 5: Mesh sizes at which a discretization error of $10^{-3}$ is attained for $\mathbb{E}(k_{\text{eff},h})$, as well as numbers of QMC or MC samples and CPU-time (2GHz Intel T7300) necessary to obtain a 95% confidence level that is also $< 10^{-3}$. The last column shows the corresponding estimates for $\mathbb{E}(k_{\text{eff},h})$ in each case.

| | Estimates for $\mathbb{E}(T_h)$ | | Estimates for $\mathbb{E}(T - T_h)$ | |
|---------|----------------------|----------------------|----------------------|----------------------|
| $1/h$ | Case 2 | Case 5 | Case 2 | Case 5 |
| 33 | $1.307184 \pm 3.7[\text{-}4]$ | $1.572801 \pm 9.8[\text{-}4]$ | $0.004392 \pm 5.1[\text{-}4]$ | $0.065798 \pm 1.6[\text{-}3]$ |
| 65 | $1.304328 \pm 5.1[\text{-}4]$ | $1.533775 \pm 1.3[\text{-}3]$ | $0.001536 \pm 6.2[\text{-}4]$ | $0.026772 \pm 1.8[\text{-}3]$ |
| 129 | $1.303665 \pm 3.7[\text{-}4]$ | $1.516761 \pm 1.6[\text{-}3]$ | $0.000873 \pm 5.1[\text{-}4]$ | $0.009758 \pm 2.0[\text{-}3]$ |
| 257 | $1.303130 \pm 3.2[\text{-}4]$ | $1.511971 \pm 1.2[\text{-}3]$ | $0.000338 \pm 4.8[\text{-}4]$ | $0.004968 \pm 1.7[\text{-}3]$ |
| $\infty$ | $1.302792 \pm 3.6[\text{-}4]$ | $1.507003 \pm 1.2[\text{-}3]$ | | |

Table 6: Estimates for $\mathbb{E}(T_h)$ computed with QMC (with $N \approx 2.1$ million). The confidence intervals are estimated using 16 random shifts. The values in the last row are estimates for the exact value of $\mathbb{E}(T)$. The values in the last two columns are estimates of $\mathbb{E}(T - T_h)$, each obtained via linear regression (with $\beta = 1.35$).

Thus, to get a discretization error of less than $2 \times 10^{-3}$ in each case, mesh sizes of $h = 1/65$ and of $h = 1/513$ are necessary, respectively. The numbers of samples necessary for the 95% confidence level of $\mathbb{E}(T_h)$ to be also less than $2 \times 10^{-3}$ in Case 2 are $N_{\text{QMC}} = 167500$ for QMC and $N_{\text{MC}} = 547500$ for MC. In Case 5, $N_{\text{QMC}} = 1.3$ million and $N_{\text{MC}} = 2.2$ million are necessary. As noted above, the advantage of QMC is less pronounced for this quantity of interest, but we still obtain a speedup of about 5 and 1.7 in Cases 2 and 5, respectively. The speedup is more pronounced if we require a higher accuracy. The FE error is less than $10^{-3}$ for $m = 129$ in Case 2. To obtain a similar accuracy with QMC and MC, $N_{\text{QMC}} = 246000$ and $N_{\text{MC}} = 3260000$ samples are necessary, and so QMC is about 12 times faster than MC in that case (7h versus 90h). To obtain a similar accuracy in Case 5 would require about 2 million samples on a grid with $h \approx 10^{-3}$ and about 6000 hours (that is 250 days) of processor time.

## 7.5 2-norm results

We finish by giving some results with the 2-norm covariance function, i.e., choosing $p = 2$ in (2.5) instead of $p = 1$. The convergence behavior of QMC is similar to the 1-norm case and we only present the pressure head in Case 2 and the effective permeability in Case 5 in Figure 7. However, interestingly the FE error behaves differently with a typical convergence rate of $\beta < 1$. This may be down to the fact that the coefficient function is not grid aligned anymore in the case of the 2-norm covariance. Table 7 shows the behavior of the FE error in the expected value of the effective permeability for Cases 2 and 5. The behavior for the FE error in the breakthrough time is similar. Note that padding is required for Case 2, but not required for Case 5.

We see from Table 7 that the accuracy obtained for a fixed grid size is significantly poorer in the case of the 2-norm covariance compared to the 1-norm covariance. To obtain an FE error of less than $10^{-2}$, mesh sizes of $h = 1/65$ and $1/513$ are necessary respectively in Cases 2
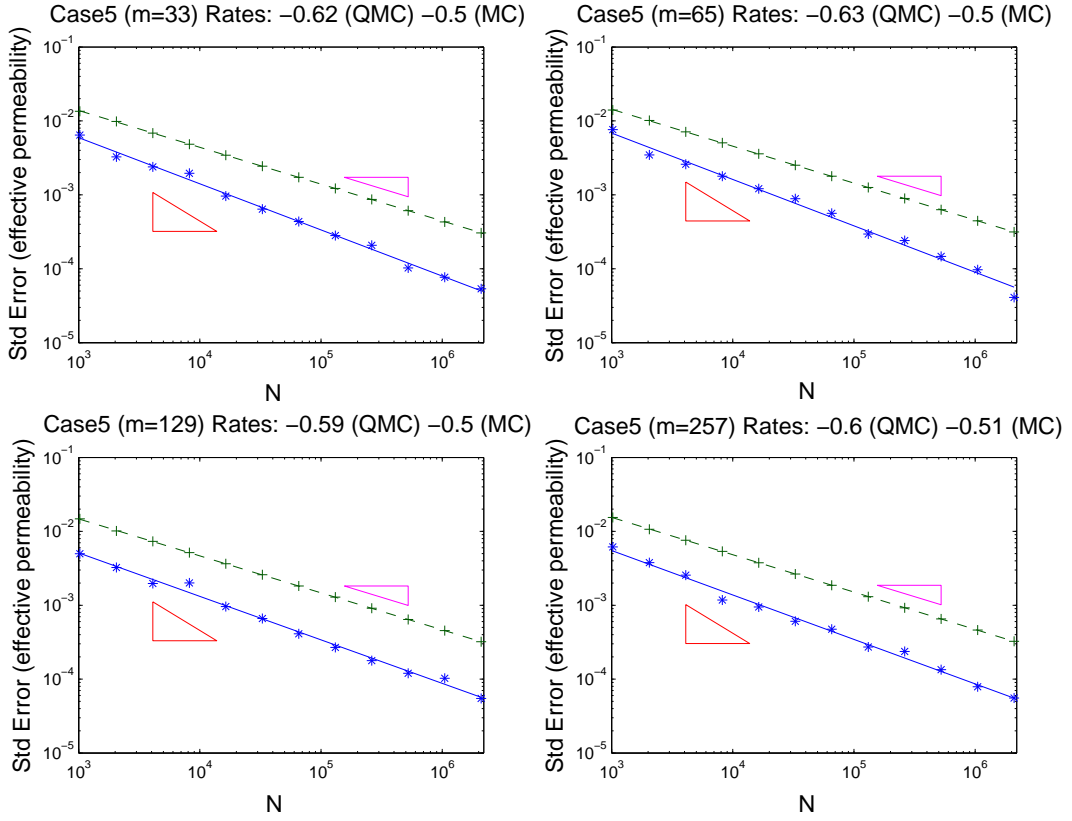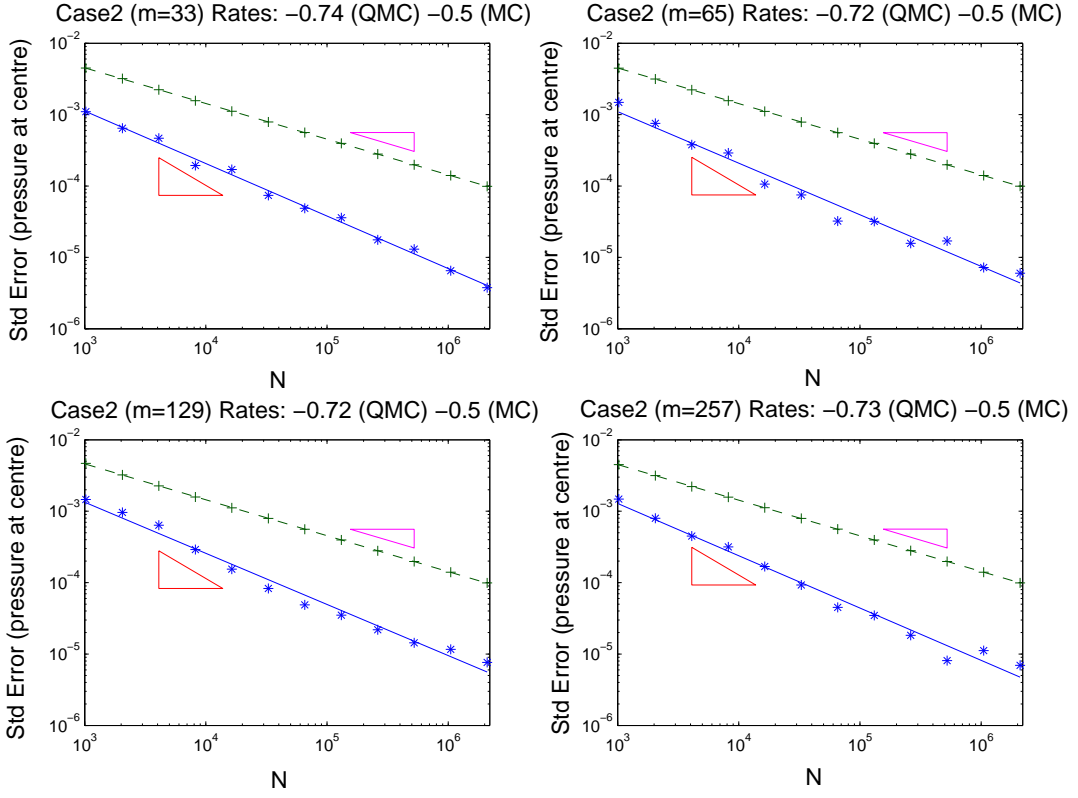
Figure 7: Standard errors $s_N$ and $s_{N,16}$ of expected pressure at centre in Case 2 (above) and of the effective permeability $k_{\text{eff},h}$ in Case 5 (below), for various values of $m$ and for the 2-norm covariance function (i.e., $p = 2$ in (2.5)).

| 1/h | Estimates for $\mathbb{E}(k_{\mathrm{eff},h})$ | | Estimates for $\mathbb{E}(k_{\mathrm{eff}} - k_{\mathrm{eff},h})$ | |
|---|---|---|---|---|
| | Case 2 | Case 5 | Case 2 | Case 5 |
| 33 | $1.118664 \pm 3.5[\text{-}5]$ | $1.001956 \pm 1.1[\text{-}4]$ | $0.015769 \pm 3.2[\text{-}4]$ | $0.090676 \pm 2.9[\text{-}3]$ |
| 65 | $1.125159 \pm 2.7[\text{-}5]$ | $1.035609 \pm 8.0[\text{-}5]$ | $0.009273 \pm 3.2[\text{-}4]$ | $0.057023 \pm 2.9[\text{-}3]$ |
| 129 | $1.128976 \pm 5.2[\text{-}5]$ | $1.059344 \pm 1.1[\text{-}4]$ | $0.005456 \pm 3.2[\text{-}4]$ | $0.033288 \pm 2.9[\text{-}3]$ |
| 257 | $1.130983 \pm 6.9[\text{-}5]$ | $1.073948 \pm 1.1[\text{-}4]$ | $0.003449 \pm 3.2[\text{-}4]$ | $0.018684 \pm 2.9[\text{-}3]$ |
| $\infty$ | $1.134433 \pm 3.2[\text{-}4]$ | $1.092632 \pm 2.9[\text{-}3]$ | | |

Table 7: Estimates for $\mathbb{E}(k_{\mathrm{eff},h})$ computed with QMC (with $N \approx 2.1 \times 10^6$) for the 2-norm covariance function (i.e., $p = 2$ in (2.5)). The confidence intervals are estimated using 16 random shifts. The values in the last row are estimates for the exact value of $\mathbb{E}(k_{\mathrm{eff}})$. The values in the last two columns are estimates of $\mathbb{E}(k_{\mathrm{eff}} - k_{\mathrm{eff},h})$ obtained via linear regression (with $\beta = 0.75$).

and 5. The numbers of samples necessary to achieve a similar accuracy for the 95% confidence level of the QMC and of the MC results are $N_{\mathrm{QMC}} = 1040$ and $N_{\mathrm{MC}} = 13900$ in Case 2, i.e., a factor of about 14. In Case 5 the corresponding numbers of samples are $N_{\mathrm{QMC}} = 1560$ and $N_{\mathrm{MC}} = 8420$. If we use $h = 1/513$ for Case 2 instead, we get a FE error of about $2 \times 10^{-3}$, and $N_{\mathrm{QMC}} = 12500$ and $N_{\mathrm{MC}} = 365000$ samples are necessary to achieve a similar accuracy in the quadrature, respectively (i.e., about 30 times less for QMC). So as before the advantage of QMC is more pronounced, if we require a higher accuracy.

# References

[1] R.J. Adler, *The Geometry of Random Fields*, Wiley, London, 1981.

[2] I. Babuška, F. Nobile and R. Tempone, A stochastic collocation method for elliptic partial differential equations with random input data, *SIAM J. Numer. Anal.* **45**, 1005–1034 (2007).

[3] I. Babuška, R. Tempone and G.E. Zouraris, Galerkin finite element approximations of stochastic elliptic partial differential equations, *SIAM J. Numer. Anal.* **42**, 800–825 (2004).

[4] F. Brezzi and M. Fortin, *Mixed and Hybrid Finite Element Methods*, Springer, New York, 1991.

[5] G. Chan and A.T.A. Wood, Algorithm AS 312: An Algorithm for simulating stationary Gaussian random fields, *Applied Statistics* **46**, 171–181 (1997).

[6] R.H. Chan and G. Strang, The asymptotic Toeplitz-circulant eigenvalue problem, *Research Report LIDS-P-1671*, Laboratory for Information and Decision Systems, M.I.T., 1987 (unpublished).

[7] K.A. Cliffe, I.G. Graham, R. Scheichl and L. Stals, Parallel computation of flow in heterogeneous media using mixed finite elements, *J. Comp. Physics* **164**, 258–282 (2000).

[8] K.A. Cliffe, I.G. Graham, R. Scheichl and L. Stals, Parallel computation of flow in heterogeneous media using mixed finite elements, *Mathematics Preprint number 99/16*, University of Bath, 1999. Available at http://www.maths.bath.ac.uk/MATHEMATICS/preprints.html .

[9] R. Cools, F.Y. Kuo, and D. Nuyens, Constructing embedded lattice rules for multivariate integration, *SIAM J. Sci. Comput.* **28**, 2162–2188 (2006).

[10] G. Dagan Solute transport in heterogeneous porous formations, *J. Fluid Mech.* **145**, 151–177 (1984).

[11] C.R. Dietrich and G.H. Newsam, Fast and exact simulation of stationary Gaussian processes through circulant embedding of the covariance matrix, *SIAM J. Sci. Comput.* **18**, 1088–1107 (1997).

[12] R.E. Ewing and J. Wang, Analysis of the Schwarz algorithm for mixed finite element methods, *RAIRO Model. Math. Anal. Num.* **26**(6), 739–756 (1992).

[13] R. Freeze, A Stochastic Conceptual Analysis of One-Dimensional Groundwater Flow in Nonuniform Homogeneous Media, *Water Resour. Res.* **11**, 725–741 (1975).

[14] M. Frigo and S.G. Johnson, The design and implementation of FFTW3, *Proceedings of the IEEE* **93** (2), 216–231 (2005).

[15] R.G. Ghanem and P.D. Spanos, *Stochastic Finite Elements*, Dover, 1991.

[16] I.G. Graham, P. Lechner and R. Scheichl. Domain decomposition for multiscale PDEs, *Numer. Math.* **106**, 589–626 (2007).

[17] F.J. Hickernell and H.S. Hong, Quasi-Monte Carlo methods and their randomisations, in R. Chan, Y.-K. Kwok, D. Yao, and Q. Zhang (eds.), *Applied Probability*, AMS/IP Studies in Advanced Mathematics, vol. 26, pp. 59–77, American Mathematical Society: Providence, 2002.

[18] F.J. Hickernell and H. Woźniakowski, Integration and approximation in arbitrary dimensions, *Adv. Comput. Math.* **12**, 25–58 (2000).

[19] N.J. Higham, *Accuracy and Stability of Numerical Algorithms*, SIAM Philadelphia, 2002.

[20] R.J. Hoeksema and P.K. Kitanidis, Analysis of the Spatial Structure of Properties of Selected Aquifers, *Water Resour. Res.* **21**, 536–572 (1985).

[21] S. Joe and F.Y. Kuo, Construction of Sobol′ sequences with better two-dimensional projections, *SIAM J. Sci. Comput.*, **30**, 2635–2654 (2008).

[22] F.Y. Kuo, Component-by-component constructions achieve the optimal rate of convergence, *J. Complexity* **19**, 301–320 (2003).

[23] F.Y. Kuo and I.H. Sloan, Lifting the curse of dimensionality, *Notices AMS*, 1320–1328 (2005).

[24] P. Monk, *Finite Element Methods for Maxwell's Equations*, Clarendon Press, Oxford, 2003.

[25] R.L. Naff, D.F. Haley, and E.A. Sudicky, High-resolution Monte Carlo simulation of flow and conservative transport in heterogeneous porous media 1. Methodology and flow results, *Water Resour. Res.*, **34**, 663–677 (1998).

[26] R.L. Naff, D.F. Haley, and E.A. Sudicky, High-resolution Monte Carlo simulation of flow and conservative transport in heterogeneous porous media 2. Transport Results, *Water Resour. Res.*, **34**, 679–697 (1998).

[27] H. Niederreiter *Random Number Generation and quasi-Monte Carlo methods*, SIAM, Philadelphia, 1994.

[28] Nirex 95: A Preliminary Analysis of the Groundwater Pathway for a Deep Repository at Sellafield. *Nirex Science Report S/95/012*, United Kingdom Nirex Limited, 1995.

[29] F. Nobile, R. Tempone and C.G. Webster, A sparse grid stochastic collocation method for partial differential equations with random input data, *SIAM J. Numer. Anal.* **46**, 2309–2345 (2008).

[30] F. Nobile, R. Tempone and C.G. Webster, An anisotropic sparse grid stochastic collocation method for partial differential equations with random input data, *SIAM J. Numer. Anal.* **46**, 2411–2442 (2008).

[31] D. Nuyens and R. Cools, Fast algorithms for component-by-component construction of rank-1 lattice rules in shift-invariant reproducing kernel Hilbert spaces *Math. Comput.* **75**, 903–920 (2006).

[32] C. Pechstein and R. Scheichl, Analysis of FETI Methods for Multiscale PDEs, *Numer. Math.* **111**(2), 293–333 (2008).

[33] C.E. Powell and H.C. Elman Block-diagonal preconditioning for spectral stochastic finite element systems, *IMA J. Numer. Anal.* **29**(2), 350–375 (2009).

[34] J. Ruge and K. Stüben, Efficient solution of finite difference and finite element equations by algebraic multigrid (AMG), In: *Multigrid Methods for Integral and Differential Equations*, D.J. Paddon and H. Holstein (eds.), pp. 169–212, IMA Conference Series, Clarendon Press, Oxford (1985).

[35] R. Scheichl, *Iterative Solution of Saddle-Point Problems using Divergence-free Finite Elements with Applications to Groundwater Flow*, PhD Thesis, University of Bath, 2000.

[36] R. Scheichl, Decoupling Three-dimensional Mixed Problems Using Divergence-free Finite Elements, *SIAM J. Sci. Comput.* **23**, 1752–1776 (2002).

[37] C. Schwab and R.A. Todor, Convergence rates for sparse chaos approximations of elliptic problems with stochastic coefficients, *IMA J. Numer. Anal.* **27**, 232–261 (2007).

[38] I.H. Sloan and S. Joe, *Lattice Methods for Multiple Integration*, Oxford Scientific Publications, Oxford, 1994.

[39] I.H. Sloan, F.K. Kuo and S. Joe, Constructing randomly shifted lattice rules in weighted Sobolev spaces, *SIAM J. Numer. Anal.* **40**, 1650–1665 (2002).

[40] I.H. Sloan and H. Woźniakowski, When are quasi-Monte Carlo algorithms efficient for high-dimensional integrals?, *J. Complexity* **14** 1–33 (1998).

[41] I.M. Sobol′, On the distribution of points in a cube and the approximate evaluation of integrals, *Zh. vȳchisl. Mat. mat. Fiz.* **7**, 784–802 (1967). English translation: *U.S.S.R. Comput. Maths. Math. Phys.* **7**, 86–112 (1967).

[42] G. Strang, *An Introduction to Applied Mathematics*, Wellesley Press, 1986.

[43] X. Wang, Strong tractability of multivariate integration using quasi-Monte Carlo algorithms, *Math. Comput.* **72**, 823–838 (2002).

[44] Y.-K. Zhang and B.-M. Seo, Numerical simulations of non-ergodic solute transport in three-dimensional heterogeneous porous media, *Stochastic Environmental Research* **18**, 205–215 (2004).