# Randomized multilevel Monte Carlo for embarrassingly parallel inference

Ajay Jasra<sup>1</sup>, Kody J. H. Law<sup>2</sup>, Alexander Tarakanov<sup>2</sup>, and Fangyuan Yu<sup>1</sup>

<sup>1</sup> Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah University of Science and Technology, Thuwal, 23955, KSA. ajay.jasra@kaust.edu.sa, fangyuan.yu@kaust.edu.sa

<sup>2</sup> Department of Mathematics, University of Manchester, Manchester, M13 9PL, UK. kodylaw@gmail.com, tarakanov517@gmail.com

Abstract. This position paper summarizes a recently developed research program focused on inference in the context of data centric science and engineering applications, and forecasts its trajectory forward over the next decade. Often one endeavours in this context to learn complex systems in order to make more informed predictions and high stakes decisions under uncertainty. Some key challenges which must be met in this context are robustness, generalizability, and interpretability. The Bayesian framework addresses these three challenges, while bringing with it a fourth, undesirable feature: it is typically far more expensive than its deterministic counterparts. In the 21st century, and increasingly over the past decade, a growing number of methods have emerged which allow one to leverage cheap low-fidelity models in order to precondition algorithms for performing inference with more expensive models and make Bayesian inference tractable in the context of high-dimensional and expensive models. Notable examples are multilevel Monte Carlo (MLMC), multi-index Monte Carlo (MIMC), and their randomized counterparts (rMLMC), which are able to provably achieve a dimension-independent (including  $\infty$ -dimension) canonical complexity rate with respect to mean squared error (MSE) of 1/MSE. Some parallelizability is typically lost in an inference context, but recently this has been largely recovered via novel double randomization approaches. Such an approach delivers independent and identically distributed samples of quantities of interest which are unbiased with respect to the *infinite* resolution target distribution. Over the coming decade, this family of algorithms has the potential to transform data centric science and engineering, as well as classical machine learning applications such as deep learning, by scaling up and scaling out fully Bayesian inference.

**Keywords:** Randomization Methods; Markov chain Monte Carlo; Bayesian Inference

### 1 Introduction

The Bayesian framework begins with a statistical model characterizing the causal relationship between various variables, parameters, and observations. A canoni-

cal example in the context of inverse problems is

$$y \sim N(G_{\theta}(u), \Gamma_{\theta}), \quad u \sim N(m_{\theta}, C_{\theta}), \quad \theta \sim \pi_0,$$

where N(m, C) denotes a Gaussian random variable with mean m and covariance  $C, G_{\theta} : U \to \mathbb{R}^{m}$  is the (typically nonlinear) parameter-to-observation map,  $\theta \in \mathbb{R}^{p}$  is a vector of parameters with  $\pi_{0}$  some distribution, and the data is given in the form of *observations* y [53,54]. Nothing precludes the case where U is a function space, e.g. leading to a Gaussian process prior above, but to avoid unnecessary technicalities, assume  $U = \mathbb{R}^{d}$ . The objective is to *condition* the prior knowledge about  $(u, \theta)$  with the observed data y and recover a *posterior distribution* 

$$p(u,\theta|y) = \frac{p(u,\theta,y)}{p(y)} = \frac{p(y|u,\theta)p(u|\theta)p(\theta)}{\int_{U \times \mathbb{R}^p} p(y|u,\theta)p(u|\theta)p(\theta)dud\theta} \,.$$

Often in the context above one may settle for a slightly simpler goal of identifying a point estimate  $\theta^*$ , e.g.  $\theta^* = \operatorname{argmax}_{\theta} p(\theta|y)$  (which we note may require an intractable integration over U) and targeting  $p(u|y, \theta^*)$  instead.

In the context described above, one often only has access to an *approximation* of the map  $G_{\theta}$ , and potentially an approximation of the domain U, which may in principle be infinite dimensional. One example is the numerical solution of a system of differential equations. Other notable examples include surrogate models arising from reduced-physics or machine-learning-type approximations [47] or deep feedforward neural networks [45]. For the sake of concreteness the reader can keep this model in mind, however it is noted that the framework is much more general, for example the parameters  $\theta$  can encode the causal relationship between latent variables via a graphical model such as a deep belief network or deep Boltzmann machine [5,43].

A concise statement of the general problem of Bayesian inference is that it requires exploration of a posterior distribution  $\Pi$  from which one cannot obtain independent and identically distributed (i.i.d.) samples. Specifically, the aim is to compute quantities such as

$$\Pi_{\theta}(\varphi) := \int_{U} \varphi(u) \Pi_{\theta}(du) \,, \quad \varphi : U \to \mathbb{R} \,, \tag{1}$$

where  $\Pi_{\theta}(du) = \pi_{\theta}(u)\nu_{\theta}(du)$ ,  $\nu_{\theta}(du)$  is either Lebesgue measure  $\nu_{\theta}(du) = du$ , or one can simulate from it,  $\pi_{\theta}(u) = \gamma_{\theta}(u)/\nu_{\theta}(\gamma_{\theta})$ , and given u one can evaluate  $\gamma_{\theta}(u)$  (or at least a non-negative unbiased estimator). Markov chain Monte Carlo (MCMC) and sequential Monte Carlo (SMC) samplers can be used for this [50]. Considering the example above with  $U = \mathbb{R}^d$ , we may take  $\nu(du) = du$  and then

$$\gamma_{\theta}(u) = |\Gamma_{\theta}|^{-1/2} |C_{\theta}|^{-1/2} \exp\left(-\frac{1}{2} |\Gamma_{\theta}^{-1/2}(y - G_{\theta}(u))|^2 - \frac{1}{2} |C_{\theta}^{-1/2}(u - m_{\theta})|^2\right), \quad (2)$$

where |A| denotes the determinant for a matrix  $A \in \mathbb{R}^n$ . Note we have used a subscript for  $\theta$ , as is typical in the statistics literature to denote that everything

3

is conditional on  $\theta$ , and note that the  $\theta$ -dependent constants are not necessary here, per se, but it is customary to define the un-normalized target as the joint on (u, y), such that  $Z_{\theta} := \nu_{\theta}(\gamma_{\theta}) = p(y|\theta)$ . Also note that in (2), u would be referred to as a *latent variable* in the statistics and machine learning literature, and so this setup corresponds to a complex physics-informed (via  $G_{\theta}$ ) unsupervised *learning model.* Labelled data problems like regression and classification [44], as well as semi-supervised learning [40,57], can also be naturally cast in a Bayesian framework. In fact, if  $G_{\theta}(u)$  is point-wise evaluation of u, i.e.  $G_{\theta}^{i}(u) = u(x^{i})$ , for inputs or covariates  $x^i$  associated to labels  $y^i$ , and one allows U to be an infinitedimensional (reproducing kernel) Hilbert space, then standard Gaussian process (GP) regression has this form. In infinite-dimensions there is no Lebesgue density, so (2) does not make sense, but the marginal likelihood and posterior can both be computed in closed form thanks to the properties of GP [48]. Alternatively, if u are the parameters of a deep feedforward neural network [45]  $f_{\theta}(\cdot; u)$ , and  $G_{\theta}^{i}(u) = f_{\theta}(x^{i}; u)$  with Gaussian prior on u, then one has a standard Bayesian neural network model [45,5].

## 1.1 The sweet and the bitter of Bayes

Three challenges which are elegantly handled in a Bayesian framework are (a) robustness, (b) generalizability, and (c) interpretability [1,52]. Uncertainty quantification (UQ) has been a topic of great interest in science and engineering applications over the past decades, due to its ability to provide a more robust model [17,1]. A model which can extrapolate outside training data coverage is referred to as generalizable. Notice that via prior knowledge (1) and the physical model, (2) has this integrated capability by design. Interpretability is the most heavily loaded word among the three desiderata. Our definition is that the model (i) can be easily understood by the user [8], (ii) incorporates all data and domain knowledge available in a principled way [8,27], and (iii) enables inference of causal relationships between latent and observed variables [46]. The natural question is then, "Why in the age of data doesn't everybody adopt Bayesian inference for all their learning requirements?"

The major hurdle to widespread adoption of a fully Bayesian treatment of learning is the computational cost. Except for very special cases, such as GP regression [48], the solution cannot be obtained in closed form. Point estimates, Laplace approximations [51], and variational methods [38,6] have therefore taken center stage, as they can yield acceptable results very quickly in many cases. In particular, for a strongly convex objective function, gradient descent achieves exponential convergence to a local minimizer, i.e. MSE  $\propto \exp(-N)$  in N steps. Such point estimates are still suboptimal from a Bayesian perspective, as they lack UQ. In terms of computation of (1), Monte Carlo (MC) methods are able to achieve exact inference in (1) in general [41,50]. In the case of i.i.d. sampling, MC methods achieve the canonical, dimension-independent convergence rate of MSE  $\propto 1/N$ , for N-sample approximations, without any smoothness assumptions and out-of-the-box<sup>3</sup>. Quadrature methods [16] and quasi-MC [9] are able to achieve improvements over MC rates, however the rates depend on the dimension and the smoothness of the integrand.

A curse of dimensionality can still hamper application of MC methods through the constant and the cost of simulation, meaning it is rare to achieve canonical *complexity* of cost  $\propto 1/\text{MSE}$  for non-trivial applications. Usually this is manifested in the form of a penalty in the exponent, so that cost  $\propto \text{MSE}^{-a}$ , for a > 2. A notable exception is MLMC [23,19] and MIMC [22] methods, and their randomized counterparts rMLMC [49,55] and rMIMC [12], which are able to achieve dimension-independent canonical complexity for a range of applications. These estimators are constructed by using a natural telescopic sum identity and constructing coupled increment estimators of decreasing variance. As an added bonus, the randomized versions eliminate discretization bias *entirely*, and deliver estimates with respect to the limiting *infinite-resolution distribution*.

In the context of inference problems, i.i.d. sampling is typically not possible and one must resort to MCMC or SMC [50]. This makes application of (r)MLMC and (r)MIMC more complex. Over the past decade, there has been an explosion of interest in applying these methods to inference, e.g. see [25,15,3,26,33] for examples of MLMC and [31,35] for MIMC. A notable benefit of MC methods is easy *parallelizability*, however typically MLMC and MIMC methods for inference are much more synchronous, or even serial in the case of MCMC. A family of rMLMC methods have recently been introduced for inference [32,36,24], which largely recover this lost parallelizability, and deliver i.i.d. samples that are unbiased with respect to the limiting infinite resolution target distribution *in the inference context*. In other words, the expectation of the resulting estimators are free from any approximation error. The first instance of rMLMC for inference was [10], and the context was different to the above work – in particular, consistent estimators are constructed that are free from discretization bias.

The rest of this paper is focused on these novel parallel rMLMC methods for inference, which are able to achieve the gold standard of Bayesian posterior inference with canonical complexity rate 1/MSE. In the age of data and increasing parallelism of supercomputer architecture, these methods are prime candidates to become a staple, if not the defacto standard, for inference in data-centric science and engineering applications. Section 2 describes some technical details of the methods, Section 3 presents a specific motivating example Bayesian inverse problem and some compelling numerical results, and Section 4 concludes with a call to action and roadmap forward for this exciting research program.

# 2 Technical Details of the methodology

The technical details of the methodology will be sketched in this section. The idea is to give an accessible overview and invitation to this exciting methodol-

<sup>&</sup>lt;sup>3</sup> This is the same rate achieved by gradient descent for general non-convex smooth objective functions. In fact, the success of deep neural networks for learning high-dimensional functions has been attributed to this dimension-independence in [56].

ogy. The interested reader can find details in the references cited. With respect to the previous section, the notation for  $\theta$  will be suppressed – the concerned reader should imagine either everything is conditioned on  $\theta$  or it has been absorbed into  $u \leftarrow (u, \theta)$ . Subsection 2.1 sketches the MLMC idea, and some of the challenges, strategies for overcoming them, and opportunities in the context of inference. Subsection 2.2 sketches the rMLMC idea, and some of the challenges, strategies for overcoming them, and opportunities in the context of inference. Finally subsection 2.3 briefly sketches MIMC.

### 2.1 Multilevel Monte Carlo

As mentioned above, for problems requiring approximation, MLMC methods are able to achieve a *huge speedup* in comparison to the naive approach of using a single fixed approximation, and indeed in some cases canonical complexity of  $\cot \infty 1/MSE$ . These methods leverage a range of successive approximations of increasing cost and accuracy. In a simplified description, most MLMC theoretical results rely on underlying assumptions of

- (i) a hierarchy of targets  $\Pi_l$ ,  $l \ge 0$ , of increasing cost, such that  $\Pi_l \to \Pi$  as  $l \to \infty$ ;
- (ii) a coupling  $\Pi^l$  s.t.  $\forall A \subset U$ ,

$$\int_{A\times U}\Pi^l(du,du')=\Pi_l(A)\,,\quad\text{and}\quad\int_{U\times A}\Pi^l(du,du')=\Pi_{l-1}(A);$$

(iii) the coupling is such that

$$\int |\varphi(u) - \varphi(u')|^2 \Pi^l(du, du') \le Ch_l^\beta, \qquad (3)$$

and the cost to simulate from  $\Pi^l$  is proportional to  $Ch_l^{-\zeta}$ , for some  $h_l > 0$ s.t.  $h_l \to 0$  as  $l \to \infty$ , and  $C, \beta, \zeta > 0$  independent of l.

Now one leverages the telescopic sum

$$\Pi(\varphi) = \underbrace{\sum_{l=0}^{L} \Delta_l(\varphi)}_{\text{approximation}} + \underbrace{\sum_{l=L+1}^{\infty} \Delta_l(\varphi)}_{\text{bias}}, \qquad (4)$$

where  $\Delta_l(\varphi) = \Pi_l(\varphi) - \Pi_{l-1}(\varphi)$ ,  $\Pi_{-1} \equiv 0$ , by approximating the first term,  $\Pi_L(\varphi)$ , using i.i.d. samples from the couplings  $\Pi^l$ ,  $l = 0, \ldots, L$ . The second term is the bias= $\Pi(\varphi) - \Pi_L(\varphi)$ . This allows one to optimally balance cost with more samples on coarse/cheap levels, and a decreasing number of samples as lincreases, to construct a multilevel estimator  $\widehat{\Pi}(\varphi)$  that achieves a given mean square error (MSE),

$$\mathbb{E}(\widehat{\Pi}(\varphi) - \Pi(\varphi))^2 = \text{variance} + \text{bias}^2,$$



(a) Few high fidelity (high-cost) simulations are combined with many at low-fidelity (low cost).

(b) MSE vs Cost: MLSMC vs SMC for elliptic PDE, illustrating the large gain in efficiency (smaller Cost for a given MSE). [3]

Fig. 1. Synopsis of MLMC methods: a family of models, including coarse-resolution approximation of differential equations, surrogates, etc. (a) can be combined in the MLMC framework to yield improved complexity cost  $\propto 1/MSE$  (b).

more efficiently than a single level method. A schematic is given in Fig. 1(a).

The MLMC estimator is defined as

$$\widehat{Y} = \sum_{l=0}^{L} \frac{1}{N_l} \sum_{i=1}^{N_l} Y_l^i \,, \tag{5}$$

where  $Y_l^i = \varphi(U_l^i) - \varphi(U_{l-1}^i)$  and  $(U_l, U_{l-1})^i \sim \Pi^l$  for  $l \geq 1$ ,  $Y_0^i = \varphi(U_0^i)$ ,  $U_0^i \sim \Pi_0$ , and L and  $\{N_l\}_{l=0}^L$  are chosen to balance the bias and variance. In particular,  $L \propto \log(\text{MSE})$  and  $N_l \propto h_l^{(\beta+\zeta)/2}$ . In the canonical regime where  $\beta > \zeta$  one achieves the canonical complexity of cost  $\propto 1/\text{MSE}$ . If  $\beta \leq \zeta$ , there are penalties. See [19] for details. Note that for the theory above, controlling the bias requires only  $\alpha > 0$  such that

$$\left|\int \varphi(u) - \varphi(u')\Pi^l(du, du')\right| \le Ch_l^{\alpha},$$

however it is clear that Jensen's inequality provides  $\alpha \geq \beta/2$ , which is suitable for the purposes of this exposition. There are notable exceptions where one can achieve  $\alpha > \beta/2$ , e.g. Euler-Maruyama or Milstein simulation of SDE [18], and this of course provides tighter results.

Note that the assumptions above can be relaxed substantially if one sacrifices a clean theory. In particular, the models  $\Pi_l$  need not be defined hierarchically in terms of a small parameter  $h_l$  corresponding to "resolution", as long as  $h_l^\beta$ and  $h_l^{-\zeta}$  in assumption (iii) above can be replaced with  $V_l$  and  $C_l$ , respectively, such that  $V_l \to 0$  as  $C_l \to \infty$  in some fashion. Indeed in practice one need not ever consider the limit and can work with a finite set of models within the same framework, as is advocated in the related multifidelity literature (see e.g. [47]). **MLMC for inference.** In the context of inference, it is rare that one can achieve i.i.d. samples from couplings  $\Pi^l$ . As described in Section 1, one more often only has access to (unbiased estimates of) the un-normalized target and must resort to MCMC or SMC. In the canonical regime  $\beta > \zeta$  the theory can proceed in a similar fashion provided one can obtain estimators  $\hat{Y}_l^N$  such that for some  $C, \beta > 0$  and q = 1, 2

$$\mathbb{E}\left[\hat{Y}_{l}^{N} - \left(\Pi_{l}(\varphi) - \Pi_{l-1}(\varphi)\right)\right]^{q} \leq C \frac{h_{l}^{\beta q/2}}{N}.$$
(6)

In the sub-canonical regime, the situation is slightly more complex.

Achieving such estimates with efficient inverse MC methods has been the focus of a large body of work. These methods can be classified according to 3 primary strategies: importance sampling [25,3,2,42,37], coupled algorithms [15,26,29,21,34], and approximate couplings [30,31,35]. See e.g. [33] for a recent review. Importance sampling estimators are the simplest, and they proceed by expressing the desired increment in terms of expectation with respect to one of the levels. Its applicability is therefore limited to cases where the importance weights can be calculated or estimated. Coupled algorithms targeting the coarse and fine targets, respectively. These are in some sense the most natural, and in principle the most general, but it can be deceptively tricky to get them to work correctly. Approximate coupling is the most straightforward strategy and can also be quite versatile. In this case, one abandons exactness with respect to coarse and fine marginals, and aims only to achieve well-behaved weights associated to a change of measure with respect to an approximate coupling.

#### 2.2 Randomized Multilevel Monte Carlo

Randomized MLMC (rMLMC) is defined similarly to (5) except with a notable difference. Define a categorical distribution  $\mathbf{p} = (\mathbf{p}_0, \mathbf{p}_1, ...)$  on  $\mathbb{Z}_+$  and let  $L^i \sim \mathbf{p}$ , and  $Y_{L^i}^i$  as above. The single term estimator [49] is defined as

$$Z^i = \frac{Y^i_{L^i}}{\mathbf{p}_{L^i}} \,. \tag{7}$$

Notice that, as a result of (4),

$$\mathbb{E}Z^{i} = \sum_{l=0}^{\infty} \mathbf{p}_{l} \mathbb{E}\left(\frac{Y_{l}^{i}}{\mathbf{p}_{l}}\right) = \Pi_{0}(\varphi) + \sum_{l=1}^{\infty} \Pi_{l}(\varphi) - \Pi_{l-1}(\varphi) = \Pi(\varphi),$$

i.e. this estimator is free from discretization bias. The corresponding rMLMC estimator is given by

$$\widehat{Z} = \frac{1}{N} \sum_{i=1}^{N} Z^{i} = \sum_{l=0}^{\infty} \frac{1}{N \mathbf{p}_{l}} \sum_{i;L_{i}=l} Y_{l}^{i} \,.$$
(8)

It is easy to see that  $\mathbb{E}\#\{i; L_i = l\} = N\mathbf{p}_l$  and  $\#\{i; L_i = l\} \to N\mathbf{p}_l$  as  $N \to \infty$ , and the optimal choice level of distribution is analogous to level selection above,  $\mathbf{p}_l \propto N_l$ , with  $N_l$  as in (5). Despite the infinite sum above, this estimator does not incur infinite cost for finite N, because only finitely many summands are nonzero. Furthermore,  $\mathbf{p}_l \to 0$ , so higher levels are simulated rarely and the expected cost is also typically finite. See [49] for further details and other variants.

**rMLMC for inference** In the inference context, one typically does not have access to *unbiased* estimators of  $Y_{L^i}$ , and rather  $\mathbb{E}(\hat{Y}_l^N) \neq \Pi_l(\varphi) - \Pi_{l-1}(\varphi)$ . In the finite *L* case, one can get away with this provided (6) holds, however rMLMC methods rely on this property. In the work [10], SMC is used to construct unbiased estimators of increments with respect to the *un-normalized* target (a well-known yet rather remarkable feature of SMC methods [14]), and subsequently a ratio estimator is used for posterior expectations, which are hence biased (for finite *N*) but consistent (in the limit  $N \to \infty$ ) with respect to the infinite-resolution ( $L = \infty$ ) target. Subsequently it has been observed that another inner application of the methodology presented above in Section 2.2 allows one to *transform a consistent estimator into an unbiased estimator* [36,32].

In particular, suppose one can couple two estimators  $\hat{Y}_l^N$  and  $\hat{Y}_l^{N'}$ , with N' > N, that marginally satisfy (6), and such that the resulting estimator satisfies, for q = 1, 2,

$$\mathbb{E}\left[\hat{Y}_l^N - \hat{Y}_l^{N'}\right]^q \le C \frac{h_l^{\beta q/2}}{N} \,. \tag{9}$$

Introduce inner levels  $N_k$ ,  $k \geq 1$ , such that  $N_k \to \infty$  as  $k \to \infty$ , and another categorical distribution  $\mathbf{p} = (\mathbf{p}_0, \mathbf{p}_1, \dots)$  on  $\mathbb{Z}_+$ . Now let  $K^i \sim \mathbf{p}$ ,  $L^i \sim \mathbf{p}$  and simulate  $\hat{Y}_{L^i}^{N_{K^i}}, \hat{Y}_{L^i}^{N_{K^i-1}}$  as above. The resulting *doubly-randomized* single term estimator is given by

$$Z^{i} = \frac{1}{\mathbf{p}_{L^{i}}\mathbf{p}_{K^{i}}} \left(\hat{Y}_{L^{i}}^{N_{K^{i}}} - \hat{Y}_{L^{i}}^{N_{K^{i}-1}}\right).$$
(10)

Now, as above,

$$\mathbb{E}\left[\frac{1}{\mathsf{p}_{K^{i}}}\left(\hat{Y}_{l}^{N_{K^{i}}}-\hat{Y}_{l}^{N_{K^{i}-1}}\right)\right]=\Pi_{l}(\varphi)-\Pi_{l-1}(\varphi),$$

and hence  $\mathbb{E}Z^i = \Pi(\varphi)$ . Furthermore, the estimators (10) can be simulated i.i.d. In other words, the embarrassingly parallel nature of classical MC estimators is restored, as well as all the classical results relating to i.i.d. random variables, such as the central limit theorem.

The work [36] leverages such a doubly randomized estimator for online particle filtering in the framework of [29]. The work [32] uses a so-called *coupled* sum variant in the framework of MLSMC samplers [3]. Both of these estimators suffer from the standard limiting MC convergence rate with respect to the inner randomization, which is sub-canonical. In other words the cost to achieve an estimator at level K is  $\mathcal{O}(N_K)$  and the error is  $\mathcal{O}(N_K^{-1})$ . As a result, it is not possible to achieve finite variance and finite cost, and one must settle for finite variance and finite cost with high probability [49]. In practice, one may truncate the sum at finite  $K_{\text{max}}$  to ensure finite cost, and accept the resulting bias.

**TMLMCMC** An alternative incarnation of the inner randomization can be used in the context of MCMC, relying on the unbiased MCMC introduced in [28], which is based on the approach of [20]. In [28] one couples a pair of MCMCs  $(U_n, U'_n)$  targeting the same distribution  $\Pi$  in such a way that they (i) have the same distribution at time n,  $U_n \stackrel{\mathcal{D}}{\sim} U'_{n+1}$ , (ii) meet in finite time  $\mathbb{E}(\tau) < \infty$ ,  $\tau = \inf\{n; U_n = U'_n\}$ , and (iii) remain identical thereafter. An unbiased estimator is then obtained via

$$\widehat{X} = \varphi(U_{n^*}) + \sum_{n=n^*+1}^{\infty} \varphi(U_n) - \varphi(U'_n)$$
$$= \varphi(U_{n^*}) + \sum_{n=n^*+1}^{\tau} \varphi(U_n) - \varphi(U'_n).$$

It is clear that in expectation the sum telescopes, giving the correct expectation  $\mathbb{E}\hat{X} = \mathbb{E}(\varphi(U_{\infty})) = \Pi(\varphi)$ . Such estimators can be simulated i.i.d., which removes the fundamental serial roadblock of MCMC, and the finite meeting time ensures finite cost. Variations of the approach allow similar efficiency to a single MCMC for a single CPU implementation, i.e. without leveraging parallelization. As above, with parallel processors, the sky is the limit.

In order to apply such technology to the present context, one couples a pair of coupled chains  $(U_{n,l}, U_{n,l-1}, U'_{n,l}, U'_{n,l-1})$  such that

$$U_{n,l}, U_{n,l-1} \stackrel{\mathcal{D}}{\sim} U'_{n+1,l}, U'_{n+1,l-1},$$

yielding a foursome that is capable of delivering finite-cost unbiased estimators of  $\Pi_l(\varphi) - \Pi_{l-1}(\varphi)$ . Indeed we are also able to achieve estimates of the type in (6), and therefore (for suitable  $\beta$ ) rMLMC estimators with finite variance and finite cost. Note that only the intra-level pairs need to meet and remain faithful. Ultimately, the i.i.d. estimators have the following form. Simulate  $L^i \sim \mathbf{p}$  as described in Section 2.2, and define  $Z^i = \widehat{Y}_{L^i}^i/\mathbf{p}_{L^i}$ , where

$$\widehat{Y}_{l}^{i} = \varphi(U_{n^{*},l}) - \varphi(U_{n^{*},l-1}) + \sum_{n=n^{*}+1}^{\tau_{l}} \varphi(U_{n,l}) - \varphi(U_{n,l}') - \sum_{n=n^{*}+1}^{\tau_{l-1}} \left(\varphi(U_{n,l-1}) - \varphi(U_{n,l-1}')\right), \quad (11)$$

with  $\tau_{\ell} = \inf\{n; U_{n,\ell} = U'_{n,\ell}\}$ , for  $\ell = l, l-1$ . The final estimator is

$$\widehat{Z} = \frac{1}{N} \sum_{i=1}^{N} Z^{i} \,. \tag{12}$$

10 Jasra, Law, Yu

#### 2.3 Multi-index Monte Carlo

Recently, the hierarchical telescopic sum identity that MLMC is based upon has been viewed through the lense of sparse grids, for the case in which there are multiple continuous spatial, temporal, and/or parametric dimensions of approximation [22]. In other words, there is a hierarchy of targets  $\Pi_{\alpha}$ , where  $\alpha$  is a multi-index, such that  $\Pi_{\alpha} \to \Pi$  as  $|\alpha| \to \infty$ . Under a more complex set of assumptions, one can appeal instead to the identity

$$\Pi(\varphi) = \sum_{\alpha \in \mathcal{I}} \Delta_{\alpha}(\varphi) + \sum_{\alpha \notin \mathcal{I}} \Delta_{\alpha}(\varphi) \,, \quad \mathcal{I} \subset \mathbb{Z}^{d}_{+} \,,$$

where d-fold multi-increments  $\Delta_{\alpha}$  are used instead, i.e. letting  $e_j \in \mathbb{R}^d$  denote the  $j^{\text{th}}$  standard basis vector and  $\delta_j \Pi_{\alpha} := \Pi_{\alpha} - \Pi_{\alpha - e_j}$ , then  $\Delta_{\alpha} := \delta_d \circ \cdots \circ \delta_1 \Pi_{\alpha}$ (for any multi-index  $\alpha'$  with  $\alpha'_i < 0$  for some  $i = 1, \ldots, d, \Pi_{\alpha'} := 0$ ). The first term is approximated again using coupled samples and the second is the bias. Under suitable regularity conditions, this MIMC method *yields further huge speedup* to obtain a given level of error [19,22]. Some preliminary work in this direction has been done recently [31,35]. Forward randomized MIMC (rMIMC) has recently been done as well [12].

## 3 Motivating example

#### 3.1 Example of Problem

The following particular problem is presented as an example. This example is prototypical of a variety of inverse problems involving physical systems in which noisy/partial observations are made of the solution of an elliptic PDE and one would like to infer the diffusion coefficient. For example, the solution to the PDE v could represent pressure of a patch of land, subject to some forcing f (sources/sinks), and the diffusion coefficient  $\hat{u}(u)$  then corresponds to the subsurface permeability [54,53], a highly desirable quantity of interest in the context of oil recovery. Let  $D \subset \mathbb{R}^d$  with  $\partial D \in C^1$  convex and  $f \in L^2(D)$ . Consider the following PDE on D:

$$-\nabla \cdot (\hat{u}(u)\nabla v) = f, \quad \text{on } D, \tag{13}$$

$$v = 0, \quad \text{on } \partial D, \tag{14}$$

where the diffusion coefficient has the form

$$\hat{u}(x;u) = \bar{u} + \sum_{j=1}^{J} u_j \sigma_j \phi_j(x),$$
(15)

Define  $u = \{u_j\}_{j=1}^J$ , and the state space will be  $X = \prod_{j=1}^J [-1, 1]$ . Let  $v(\cdot; u)$  denote the weak solution of (1) for parameter value u. The prior is given by

 $u_j \sim U[-1,1]$  (the uniform distribution on [-1,1]) i.i.d. for  $j = 1, \ldots, J$ . It will be assumed that  $\phi_j \in C(D)$ ,  $\|\phi_j\|_{\infty} \leq 1$ , and there is a  $u_* > 0$  such that  $\bar{u} > \sum_{j=1}^J \sigma_j + u_*$ . Note that under the given assumptions,  $\hat{u}(u) > u_*$  uniformly in u. Hence there is a well-defined (weak) solution  $v(\cdot; u)$  that is bounded in  $L^{\infty}(D)$  and  $L^2(D)$  uniformly in u, and its gradient is also bounded in  $L^2(D)$ uniformly in u [11,13].

Define the following vector-valued function

$$G(u) = [\langle g_1, v(\cdot; u) \rangle, \dots, \langle g_m, v(\cdot; u) \rangle]^{\mathsf{T}},$$
(16)

where  $g_i \in L^2(D)$  for i = 1, ..., m. We note that pointwise evaluation is also permissible since  $u \in L^{\infty}(D)$ , i.e.  $g_i$  can be Dirac delta functions, however for simplicity we restrict the presentation to  $L^2(D)$ . It is assumed that the data take the form

$$y = G(u) + \xi, \quad \xi \sim N(0, \theta^{-1} \cdot \boldsymbol{I}_m), \quad \xi \perp u,$$
(17)

where  $\perp$  denotes independence. The unnormalized density  $\gamma_{\theta} : \mathsf{X} \to \mathbb{R}_+$  of u for fixed  $\theta > 0$  is given by

$$\gamma_{\theta}(u) = \theta^{m/2} \exp\left(-\frac{\theta}{2} \|G(u) - y\|^2\right).$$
(18)

The normalized density is

$$\eta_{\theta}(u) = \frac{\gamma_{\theta}(u)}{I_{\theta}} \,,$$

where  $I_{\theta} = \int_{\mathsf{X}} \gamma_{\theta}(u) du$ , and the quantity of interest is defined for  $u \in \mathsf{X}$  as

$$\varphi_{\theta}(u) := \nabla_{\theta} \log\left(\gamma_{\theta}(u)\right) = \frac{m}{2\theta} - \frac{1}{2} \|G(u) - y\|^2.$$
(19)

To motivation this particular objective function, notice that  $\gamma_{\theta}$  is chosen such that the marginal likelihood, or "evidence" for  $\theta$ , is given by  $p(y|\theta) = I_{\theta}$ . Therefore the MLE ( $\lambda = 0$ ) or MAP are given as minimizers of  $-\log I_{\theta} + \lambda R(\theta)$ , where  $R(\theta) = -\log p(\theta)$ . Assuming  $R(\theta)$  is known in closed form and differentiable, then a gradient descent method requires

$$\nabla_{\theta} \log I_{\theta} = \frac{1}{I_{\theta}} \int_{\mathsf{X}} \nabla_{\theta} \gamma_{\theta}(u) du = \frac{1}{I_{\theta}} \int_{\mathsf{X}} \underbrace{\nabla_{\theta} \log\left(\gamma_{\theta}(u)\right)}_{\varphi_{\theta}(u)} \gamma_{\theta}(u) du = \eta_{\theta}(\varphi_{\theta}(u)) \,. \tag{20}$$

Stochastic gradient descent requires only an unbiased estimator of  $\eta_{\theta}(\varphi_{\theta}(u))$ [39], which the presented rMLMC method delivers.

**Numerical approximation** The finite element method (FEM) is utilized for solution of (14) with piecewise multi-linear nodal basis functions. Let d = 1 and D = [0, 1] for simplicity. Note the approach is easily generalized to  $d \ge 1$  using products of such piecewise linear functions described below following standard

#### 12 Jasra, Law, Yu

FEM literature [7]. The PDE problem at resolution level l is solved using FEM with piecewise linear shape functions on a uniform mesh of width  $h_l = 2^{-l}$ , for  $l \ge 0$ . Thus, on the *l*th level the finite-element basis functions are  $\{\psi_i^l\}_{i=1}^{2^l-1}$  defined as (for  $x_i = i \cdot 2^{-l}$ ):

$$\psi_i^l(x) = \begin{cases} (1/h_l)[x - (x_i - h_l)] \text{ if } x \in [x_i - h_l, x_i], \\ (1/h_l)[x_i + h_l - x] \text{ if } x \in [x_i, x_i + h_l]. \end{cases}$$

To solve the PDE,  $v^l(x) = \sum_{i=1}^{2^l-1} v_i^l \psi_i^l(x)$  is plugged into (1), and projected onto each basis element:

$$-\left\langle \nabla \cdot \left( \hat{u} \nabla \sum_{i=1}^{2^l-1} v_i^l \psi_i^l \right), \psi_j^l \right\rangle = \langle f, \psi_j^l \rangle,$$

resulting in the following linear system:

$$\boldsymbol{A}^{l}(u)\boldsymbol{v}^{l}=\boldsymbol{f}^{l},$$

where we introduce the matrix  $\mathbf{A}^{l}(u)$  with entries  $A_{ij}^{l}(u) = \langle \hat{u} \nabla \psi_{i}^{l}, \nabla \psi_{j}^{l} \rangle$ , and vectors  $\mathbf{v}^{l}, \mathbf{f}^{l}$  with entries  $v_{i}^{l} = \langle v, \psi_{i}^{l} \rangle$  and  $f_{i}^{l} = \langle f, \psi_{i}^{l} \rangle$ , respectively.

Define  $G^l(u) = [\langle g_1, v^l(\cdot; u) \rangle, \dots, \langle g_m, v^l(\cdot; u) \rangle]^{\intercal}$ . Denote the corresponding approximated un-normalized density by

$$\gamma_{\theta}^{l}(u) = \theta^{m/2} \exp\left\{-\frac{\theta}{2} \|G^{l}(u) - y\|^{2}\right\},\tag{21}$$

and the approximated normalized density by  $\eta_{\theta}^{l}(u) = \gamma_{\theta}^{l}(u)/I_{\theta}^{l}$ , where  $I_{\theta}^{l} = \int_{\mathbf{X}} \gamma_{\theta}^{l}(u) du$ . Furthermore, define

$$\varphi_{\theta}^{l}(u) := \nabla_{\theta} \log\left(\gamma_{\theta}^{l}(u)\right) = \frac{m}{2\theta} - \frac{1}{2} \|G^{l}(u) - y\|^{2}.$$
(22)

It is well-known that under the stated assumptions  $v^l(u)$  converges to v(u) as  $l \to \infty$  in  $L^2(D)$  (as does its gradient), uniformly in u [7,11], with the rate  $h_l^{\beta/2}$ ,  $\beta = 4$ . In a forward UQ context, this immediately provides (3) for Lipschitz functions of v, with  $\beta = 4$ . Furthermore, continuity ensures  $\gamma_{\theta}^l(u)$  converges to  $\gamma_{\theta}(u)$  and  $\varphi_{\theta}^l(u)$  converges to  $\varphi_{\theta}(u)$  uniformly in u as well. See also [4,2] for further details. This allows one to achieve estimates of the type (6) in the inference context.

#### 3.2 Numerical results

This section is for illustration purposes and reproduces results from [24], specifically Section 4.1.2 and Figure 5. The problem specified in the previous section is considered with forcing f(x) = 100x. The prior specification of  $u = (u_1, u_2)$  is taken as J = 2,  $\bar{u} = 0.15$ ,  $\sigma_1 = 1/10$ ,  $\sigma_2 = 1/40$ ,  $\phi_1(t) = \sin(\pi x)$  and  $\phi_2(t) = \cos(2\pi x)$ . For this particular setting, the solution v is continuous and

hence point-wise observations are well-defined. The observation function G(x) in (16) is chosen as  $g_i(v(u)) = v(0.01+0.02(i-1); u)$  for  $i \in \{1, \ldots, m\}$  with m = 50. The FEM scheme in Section 3.1 is employed with mesh width of  $l \leftarrow l+l_0$ , where  $l_0 = 3$ . Using a discretization level of l = 10 to approximate G(x) with  $G_l(x)$ , x = (0.6, -0.4) and  $\theta = 1$ , observations  $y \in \mathbb{R}^m$  are simulated from (17).

The estimators  $\hat{Y}_{L^i}^i$  are computed using a reflection maximal coupling of pCN kernels, as described in [24]. The left panel of Figure 2 illustrates that averaging single term estimators (11) as in (12) yields a consistent estimator that converges at the canonical Monte Carlo rate of 1/MSE.

Consider now inference for  $\theta$  in the Bayesian framework, under a prior  $p(\theta)$ specified as a standard Gaussian prior on  $\log \theta$ . A stochastic gradient ascent algorithm is initialized at  $\theta^{(0)} = 0.1$  to compute the maximum a posteriori probability (MAP) estimator  $\theta_{\text{MAP}} \in \arg \max p(\theta)I_{\theta}$ , simulated by subtracting  $\nabla_{\theta}R(\theta)$  from the estimator of (20) given by  $Z^i$  defined above and in (11). The right panel of Figure 2 displays convergence of the stochastic iterates to  $\theta_{\text{MAP}}$ . An estimator following [32], of the type in (10), is also shown here, using the algorithm in [4] instead of coupled MCMC. The plot shows some gains over [32] when the same learning rates are employed.

**Parallel implementation.** An example is now presented to illustrate the parallel improvement of these methods on multiple cores. These results are borrowed from [36] for (online) filtering of partially observed diffusions. In particular, an estimator of the form (10) is constructed, in which each  $\hat{Y}_{L^i}^{N_{K^i}}$  is a coupled particle filter increment estimator at resolution  $L^i$  and with  $K^i$  particles, for  $i = 1, \ldots, N$ , and these estimators are then averaged as in (12). The parallel performance is assessed with up to  $1000(\leq N)$  MPI cores on the KAUST supercomputer Shaheen. A Python notebook that implements the unbiased estimator both on a single core and multiple cores can be found in the following Github link: https://github.com/fangyuan-ksgk/Unbiased-Particle-Filter-HPC-.

To demonstrate the parallel scaling power, various numbers of processors  $M \in \{1, 5, 10, 20, 50, 100, 500, 1000\}$  are used, with  $N = 10^3 M$ . The serial computation time to obtain the estimator on a single core is recorded, as well as the parallel computation time on M cores. The parallel speedup is defined as the ratio of cost for serial implementation and the cost for parallel implementation, and the parallel efficiency is given by the ratio of parallel speedup and the number of parallel cores M.

The results are shown in Figure 3, which shows almost perfect strong scaling for up to 1000 MPI cores, for this level of accuracy. It is important to note that there will be a limitation to the speedup possible, depending upon the accuracy level. In particular, the total simulation time is limited by the single most expensive sample required. Therefore, it will not be possible to achieve  $MSE \propto \varepsilon^2$  in  $\mathcal{O}(1)$  time, even with arbitrarily many cores.



Fig. 2. Elliptic Bayesian inverse problem of Section 3.2 Left: accuracy (minus MSE) against number of single term samples N. The samples were simulated in serial on a laptop, but can all be simulated in parallel. Right: convergence of stochastic gradient iterates  $\theta^{(n)}$  to the maximum a posteriori probability estimator  $\theta_{MAP}$ . The learning rates considered here are  $\alpha_n = \alpha_1/n$ . The red curve corresponds to the unbiased MLSMC algorithm of [32] for comparison.



Fig. 3. Parallel Speedup and Parallel Efficiency against number of MPI cores for the unbiased particle filter from [36].

16 Jasra, Law, Yu

## 4 Conclusion and path forward

This position paper advocates for the widespread adoption of Bayesian methods for performing inference, especially in the context of complex science and engineering applications, where high-stakes decisions require robustness, generalizability, and interpretability. Such methods are rapidly gaining momentum in science and engineering applications, following an explosive interest in UQ, in concert with the data deluge and emerging fourth paradigm of data-centric science and engineering. Meanwhile, in the field of machine learning and AI the value of Bayesian methods has been recognized already for several decades. There it is widely accepted that the Bayesian posterior is the gold standard, but the community has largely converged on variational approximations or even point estimators as surrogates, due to complexity limitations.

Here a family of embarrassingly parallel rMLMC simulation methods are summarized. The methods are designed for performing exact Bayesian inference in the context where only approximate models are available, which includes a wide range of problems in physics, biology, finance, machine learning, and spatial statistics. Canonical complexity is achieved. Important priorities going forward are: (i) continued development of novel instances of this powerful class of algorithms, (ii) adaptation to specific large scale application contexts across science, engineering, and AI, and (iii) automation of the methods and the design of usable software to enable deployment on a large scale and across applications in science, engineering, and AI, ideally by practitioners and without requiring an expert.

Acknowledgements. KJHL and AT were supported by The Alan Turing Institute under the EPSRC grant EP/N510129/1. AJ and FY acknowledge KAUST baseline support.

## References

- Nathan Baker, Frank Alexander, Timo Bremer, Aric Hagberg, Yannis Kevrekidis, Habib Najm, Manish Parashar, Abani Patra, James Sethian, Stefan Wild, et al. Workshop report on basic research needs for scientific machine learning: Core technologies for artificial intelligence. Technical report, USDOE Office of Science (SC), Washington, DC (United States), 2019.
- Alexandros Beskos, Ajay Jasra, Kody J. H. Law, Youssef Marzouk, and Yan Zhou. Multilevel sequential Monte Carlo with dimension-independent likelihood-informed proposals. SIAM/ASA Journal on Uncertainty Quantification, 6(2):762–786, 2018.
- Alexandros Beskos, Ajay Jasra, Kody J. H. Law, Raul Tempone, and Yan Zhou. Multilevel sequential Monte Carlo samplers. *Stochastic Processes and their Applications*, 127(5):1417–1440, 2017.
- Alexandros Beskos, Ajay Jasra, Kody J. H. Law, Raul Tempone, and Yan Zhou. Multilevel sequential Monte Carlo samplers. *Stochastic Processes and their Applications*, 127(5):1417–1440, 2017.
- 5. Christopher M Bishop. Pattern recognition and machine learning. Springer, 2006.

- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Susanne Brenner and Ridgway Scott. The mathematical theory of finite element methods, volume 15. Springer Science & Business Media, 2007.
- 8. Alan Bundy and etal. Explainable AI: the basics, 2019.
- Russel E Caflisch et al. Monte carlo and quasi-monte carlo methods. Acta numerica, 1998:1–49, 1998.
- Neil Chada, Jordan Franks, Ajay Jasra, Kody J. H. Law, and Matti Vihola. Unbiased inference for discretely observed hidden markov model diffusions. *SIAM JUQ, to appear*, 2020.
- 11. Philippe G Ciarlet. The finite element method for elliptic problems. SIAM, 2002.
- Dan Crisan, Pierre Del Moral, Jeremie Houssineau, and Ajay Jasra. Unbiased multi-index Monte Carlo. Stochastic Analysis and Applications, 36(2):257–273, 2018.
- Masoumeh Dashti and Andrew M Stuart. Uncertainty quantification and weak approximation of an elliptic inverse problem. SIAM Journal on Numerical Analysis, 49(6):2524–2542, 2011.
- 14. Pierre Del Moral. Feynman-Kac formulae. Springer, 2004.
- Tim J Dodwell, Christian Ketelsen, Robert Scheichl, and Aretha L Teckentrup. A hierarchical multilevel Markov chain Monte Carlo algorithm with applications to uncertainty quantification in subsurface flow. SIAM/ASA Journal on Uncertainty Quantification, 3(1):1075–1108, 2015.
- Thomas Gerstner and Michael Griebel. Dimension-adaptive tensor-product quadrature. Computing, 71(1):65–87, 2003.
- 17. Roger Ghanem, David Higdon, and Houman Owhadi. Handbook of uncertainty quantification, volume 6. Springer, 2017.
- Michael B Giles. Multilevel Monte Carlo path simulation. Operations research, 56(3):607–617, 2008.
- 19. Michael B Giles. Multilevel Monte Carlo methods. Acta Numer., 24:259-328, 2015.
- Peter W Glynn and Chang-han Rhee. Exact estimation for markov chain equilibrium expectations. Journal of Applied Probability, 51(A):377–389, 2014.
- Alastair Gregory, Colin J Cotter, and Sebastian Reich. Multilevel ensemble transform particle filtering. SIAM Journal on Scientific Computing, 38(3):A1317– A1338, 2016.
- 22. Abdul-Lateef Haji-Ali, Fabio Nobile, and Raúl Tempone. Multi-index Monte Carlo: when sparsity meets sampling. *Numerische Mathematik*, 132(4):767–806, 2016.
- Stefan Heinrich. Multilevel Monte Carlo methods. In International Conference on Large-Scale Scientific Computing, pages 58–67. Springer, 2001.
- 24. Jeremy Heng, Ajay Jasra, Kody J. H. Law, and Alexander Tarakanov. On unbiased estimation for discretized models. *arXiv preprint arXiv:2102.12230*, 2021.
- Viet Ha Hoang, Christoph Schwab, and Andrew M Stuart. Complexity analysis of accelerated MCMC methods for Bayesian inversion. *Inverse Problems*, 29(8):085010, 2013.
- Håkon Hoel, Kody J. H. Law, and Raúl Tempone. Multilevel ensemble Kalman filtering. SIAM Journal on Numerical Analysis, 54(3):1813–1839, 2016.
- 27. The Alan Turing Institute. The AI revolution in scientific research, 2019.
- Pierre E Jacob, John O'Leary, and Yves F Atchadé. Unbiased Markov chain Monte Carlo methods with couplings. *Journal of the Royal Statistical Society: Series B* (Statistical Methodology), 82(3):543–600, 2020.

- 18 Jasra, Law, Yu
- Ajay Jasra, Kengo Kamatani, Kody J. H. Law, and Yan Zhou. Multilevel particle filters. SIAM Journal on Numerical Analysis, 55(6):3068–3096, 2017.
- 30. Ajay Jasra, Kengo Kamatani, Kody J. H. Law, and Yan Zhou. Bayesian static parameter estimation for partially observed diffusions via multilevel Monte Carlo. *SIAM Journal on Scientific Computing*, 40(2):A887–A902, 2018.
- Ajay Jasra, Kengo Kamatani, Kody J. H. Law, and Yan Zhou. A multi-index Markov chain Monte Carlo method. *International Journal for Uncertainty Quan*tification, 8(1), 2018.
- 32. Ajay Jasra, Kody J. H. Law, and Deng Lu. Unbiased estimation of the gradient of the log-likelihood in inverse problems. *Statistics and Computing*, 31(3):1–18, 2021.
- Ajay Jasra, Kody J. H. Law, and Carina Suciu. Advanced multilevel Monte Carlo methods. *International Statistical Review*, 88(3):548–579, 2020.
- Ajay Jasra, Kody J. H. Law, and Yaxian Xu. Markov chain simulation for multilevel Monte Carlo. Foundations of Data Science, 3:27, 2021.
- Ajay Jasra, Kody J. H. Law, and Yaxian Xu. Multi-index sequential Monte Carlo methods for partially observed stochastic partial differential equations. *Interna*tional Journal for Uncertainty Quantification, 11(3), 2021.
- Ajay Jasra, Kody J. H. Law, and Fangyuan Yu. Unbiased filtering of a class of partially observed diffusions. arXiv preprint arXiv:2002.03747, 2020.
- 37. Ajay Jasra, Kody J. H. Law, and Yan Zhou. Forward and inverse uncertainty quantification using multilevel Monte Carlo algorithms for an elliptic nonlocal equation. International Journal for Uncertainty Quantification, 6(6), 2016.
- Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- 39. Harold Kushner and G George Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.
- Neil Lawrence and Michael Jordan. Semi-supervised learning via gaussian processes. Advances in neural information processing systems, 17:753–760, 2004.
- Nicholas Metropolis and Stanislaw Ulam. The Monte Carlo method. Journal of the American statistical association, 44(247):335–341, 1949.
- 42. Pierre Del Moral, Ajay Jasra, Kody J. H. Law, and Yan Zhou. Multilevel sequential Monte Carlo samplers for normalizing constants. ACM Transactions on Modeling and Computer Simulation (TOMACS), 27(3):1–22, 2017.
- 43. Kevin P Murphy. Machine learning: a probabilistic perspective. MIT press, 2012.
- R. Neal. Regression and classification using Gaussian process priors. Bayesian statistics, 6:475, 1998.
- Radford M Neal. Bayesian learning for neural networks, volume 118. Springer Science & Business Media, 2012.
- Judea Pearl et al. Causal inference in statistics: An overview. Statistics surveys, 3:96–146, 2009.
- Benjamin Peherstorfer, Karen Willcox, and Max Gunzburger. Survey of multifidelity methods in uncertainty propagation, inference, and optimization. *Siam Review*, 60(3):550–591, 2018.
- Carl Edward Rasmussen and Christopher KI Williams. Gaussian processes for machine learning. 2006.
- 49. Chang-han Rhee and Peter W Glynn. Unbiased estimation with square root convergence for SDE models. *Operations Research*, 63(5):1026–1043, 2015.
- Christian Robert and George Casella. A short history of Markov chain Monte Carlo: Subjective recollections from incomplete data. *Statistical Science*, pages 102–115, 2011.

- Håvard Rue, Sara Martino, and Nicolas Chopin. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the royal statistical society: Series b (statistical methodology)*, 71(2):319– 392, 2009.
- Rick Stevens, Valerie Taylor, Jeff Nichols, Arthur Barney Maccabe, Katherine Yelick, and David Brown. Ai for science. Technical report, Argonne National Lab.(ANL), Argonne, IL (United States), 2020.
- Andrew M Stuart. Inverse problems: a Bayesian perspective. Acta numerica, 19:451–559, 2010.
- 54. Albert Tarantola. Inverse problem theory and methods for model parameter estimation. SIAM, 2005.
- Matti Vihola. Unbiased estimators and multilevel Monte Carlo. Operations Research, 66(2):448–462, 2018.
- E Weinan, Jiequn Han, and Linfeng Zhang. Integrating machine learning with physics-based modeling. Arxiv preprint. https://arxiv.org/pdf/2006.02619.pdf, 2020.
- 57. Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using Gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pages 912–919, 2003.