

1 **ON LARGE LAG SMOOTHING FOR HIDDEN MARKOV MODELS***2 JEREMIE HOUSSINEAU[†], AJAY JASRA[†], AND SUMEETPAL S. SINGH[‡]

3 **Abstract.** In this article we consider the smoothing problem for hidden Markov models (HMM).
 4 Given a hidden Markov chain $\{X_n\}_{n \geq 0}$ and observations $\{Y_n\}_{n \geq 0}$, our objective is to compute
 5 $\mathbb{E}[\varphi(X_0, \dots, X_k) | y_0, \dots, y_n]$ for some real-valued, integrable functional φ and k fixed, $k \ll n$ and
 6 for some realisation (y_0, \dots, y_n) of (Y_0, \dots, Y_n) . We introduce a novel application of the multilevel
 7 Monte Carlo (MLMC) method with a coupling based on the Knothe-Rosenblatt rearrangement. We
 8 prove that this method can approximate the afore-mentioned quantity with a mean square error
 9 (MSE) of $\mathcal{O}(\epsilon^2)$, for arbitrary $\epsilon > 0$ with a cost of $\mathcal{O}(\epsilon^{-2})$. This is in contrast to the same direct
 10 Monte Carlo method, which requires a cost of $\mathcal{O}(n\epsilon^{-2})$ for the same MSE. The approach we suggest
 11 is, in general, not possible to implement, so the optimal transport methodology of [26, 23] is used,
 12 which directly approximates our strategy. We show that our theoretical improvements are achieved,
 13 even under approximation, in several numerical examples.

14 **Key words.** Smoothing, Multilevel Monte Carlo, Optimal Transport15 **AMS subject classifications.** 62M05, 62E17

16 **1. Introduction.** Given a hidden Markov chain $\{X_n\}_{n \geq 0}$, $X_n \in \mathbf{X} \subset \mathbb{R}^d$ and
 17 observations $\{Y_n\}_{n \geq 0}$, $Y_n \in \mathbf{Y}$, we consider a probabilistic model such that for Borel
 18 $A \in \mathcal{X}$, $\mathbb{P}(X_0 \in A) = \int_A f(x) dx$, for every $n \geq 1$, $x_{0:n-1} \in \mathbf{X}^n$

$$19 \quad (1.1) \quad \mathbb{P}(X_n \in A | x_{0:n-1}) = \int_A f(x_{n-1}, x) dx$$

20 with dx Lebesgue measure and for Borel $B \in \mathcal{Y}$ and all $n \geq 0$, $(y_{0:n-1}, x_{0:n}) \in$
 21 $\mathbf{Y}^n \times \mathbf{X}^{n+1}$

$$22 \quad (1.2) \quad \mathbb{P}(Y_n \in B | y_{0:n-1}, x_{0:n}) = \int_B g(x_n, y) dy,$$

23 where we have used the compact notation $a_{k:n} = (a_k, \dots, a_n)$ for any $k, n \geq 0$ and
 24 any sequence $(a_n)_{n \geq 0}$ with the convention that the resulting vector of objects is null
 25 if $k > n$. The model defined by (1.1) and (1.2) is termed a hidden Markov model.
 26 In this article, given $y_{0:n}$, our objective is to compute $\mathbb{E}[\varphi(X_{0:k}) | y_{0:n}]$ for some real-
 27 valued, integrable functional φ and k fixed, $k \ll n$, which we refer to as large-lag
 28 smoothing. Hidden Markov models and the smoothing problem are found in many real
 29 applications, such as finance, genetics and engineering; see e.g. [4] and the references
 30 therein.

31 The smoothing problem is notoriously challenging. Firstly, $\mathbb{E}[\varphi(X_{0:k}) | y_{0:n}]$
 32 is seldom available analytically and hence numerical methods are required.
 33 Secondly, if one wants to compute $\mathbb{E}[\varphi(X_{0:k}) | y_{0:n}]$ for several values of n ,
 34 i.e. potentially recursively, then several of the well-known methods for approximation
 35 of $\mathbb{E}[\varphi(X_{0:k}) | y_{0:n}]$ can fail. For instance the particle filter (e.g. [8] and the references

*Submitted to the editors on July 2, 2018.

Funding: All authors were supported by Singapore Ministry of Education AcRF tier 1 grant R-155-000-182-114. AJ is affiliated with the Risk Management Institute, OR and analytics cluster and the Center for Quantitative Finance at NUS.

[†]Department of Statistics & Applied Probability, National University of Singapore, Singapore, 117546, SG. (housineau.j@gmail.com, staja@nus.edu.sg).

[‡]Department of Engineering, University of Cambridge, Cambridge, CB2 1PZ and The Alan Turing Institute, UK. (sss40@cam.ac.uk)

therein) suffers from the well-known path degeneracy problem (see e.g. [19]). Despite this, several methods are available for the approximation of $\mathbb{E}[\varphi(X_{0:k})|y_{0:n}]$, such as particle Markov chain Monte Carlo [1] or the PaRIS algorithm [22], which might be considered the current state-of-the-art. The latter algorithm relies on approximating $\mathbb{E}[\varphi(X_{0:k})|y_{0:n^*}]$ for some $n^* < n$ and is then justified on the basis of using *forgetting* properties of the smoother (see e.g. [4, 7]). We will extend this notion as will be explained below.

The main approach that is followed in this paper, is to utilize the multilevel Monte Carlo method (e.g. [10, 13, 12, 15]). Traditional applications of this method are associated to discretizations of continuum problems, but we adopt the framework in a slightly non-standard way. To describe the basic idea, suppose one is interested in $\mathbb{E}_\pi[\varphi(X)]$ for π a probability, φ real-valued and bounded, but, one can only hope to approximate $\mathbb{E}_{\pi_l}[\varphi(X)]$ with π_l a probability (assumed on the same space as π), $l \in \mathbb{N}$ and in some loose sense one has π_l approaches π as l grows. Now, given π_0, \dots, π_L a sequence of increasingly more ‘precise’ probability distributions on the same space, one trivially has

$$(1.3) \quad \mathbb{E}_{\pi_L}[\varphi(X)] = \mathbb{E}_{\pi_0}[\varphi(X)] + \sum_{l=1}^L \{\mathbb{E}_{\pi_l}[\varphi(X)] - \mathbb{E}_{\pi_{l-1}}[\varphi(X)]\}.$$

The approach is now to sample *dependent* couplings of (π_l, π_{l-1}) independently for $1 \leq l \leq L$ and approximate the difference $\mathbb{E}_{\pi_l}[\varphi(X)] - \mathbb{E}_{\pi_{l-1}}[\varphi(X)]$ using Monte Carlo. The term $\mathbb{E}_{\pi_0}[\varphi(X)]$ is also approximated using Monte Carlo with i.i.d. sampling from π_0 . Then, given a ‘good enough’ coupling and a characterization of the bias, for many practical problems the cost to achieve a pre-specified MSE against i.i.d. sampling from π_L and Monte Carlo, is significantly reduced. To elaborate the effectiveness of the coupling (as discussed in [11]), the main issue is to approximate (as in eq. (1.3))

$$(1.4) \quad \mathbb{E}_{\pi_l}[\varphi(X)] - \mathbb{E}_{\pi_{l-1}}[\varphi(X)] = \mathbb{E}_{\tilde{\pi}_{l,l-1}}[\varphi(X) - \varphi(Y)]$$

where $\tilde{\pi}_{l,l-1}$ is any probability on the product space (say $\mathbb{R} \times \mathbb{R}$) of the original probability measures π_l, π_{l-1} , with for any measurable $A \subseteq \mathbb{R}$, $\int_{A \times \mathbb{R}} \tilde{\pi}_{l,l-1}(d(x, y)) = \int_A \pi_l(dx)$, $\int_{\mathbb{R} \times A} \tilde{\pi}_{l,l-1}(d(x, y)) = \int_A \pi_{l-1}(dy)$. Now, if one performs i.i.d. sampling from $\tilde{\pi}_{l,l-1}$ to approximate the R.H.S. of (1.4), the variance of this approximation (of say $N \geq 1$ samples) is upper-bounded by a term of the form

$$\frac{\|\varphi\|_{\text{Lip}}}{N} \mathbb{E}_{\tilde{\pi}_{l,l-1}}[|X - Y|^2]$$

where we assume φ is Lipschitz, $|\varphi(x) - \varphi(y)| \leq \|\varphi\|_{\text{Lip}}|x - y|$. Now, the gain of MLMC is possible if the coupling can strongly correlate X, Y . In the case above, we know that the optimal coupling is that w.r.t. squared Wasserstein distance.

We leverage the idea of MLMC where the ‘level’ l corresponds to the time parameter and L is some chosen n^* , so as to achieve a given level of bias. The main issue is then how to sample from couplings which are good enough. We show that, as elaborated on above, when $d = 1$ (the dimension of the hidden state) that using the optimal coupling, in terms of squared Wasserstein distance, can yield significant improvements over the case where one directly approximates $\mathbb{E}[\varphi(X_{0:k})|y_{0:n}]$ with Monte Carlo and i.i.d sampling from the smoother. That is, for $\epsilon > 0$ given, to achieve a mean square error of $\mathcal{O}(\epsilon^2)$, the cost is $\mathcal{O}(\epsilon^{-2})$, whereas for the ordinary Monte Carlo method the cost is $\mathcal{O}(n\epsilon^{-2})$. The same conclusion with $d > 1$ can be achieved using

73 the Knothe-Rosenblatt rearrangement. The main issue with our approach is that it
 74 cannot be implemented for most problems of practical interest. However, using the
 75 transport methodology in [26], it can be approximated. We show that in numerical
 76 examples our predicted theory is verified, even under this approximation. We also
 77 compare our method directly with PaRIS, showing substantial improvement in terms
 78 of cost for a given level of MSE. Note that the transport methodology used here differs
 79 fundamentally from the “particle flow” methods discussed in [6, 3, 14] where samples
 80 from a base probability distributions are moved using an ordinary differential equation
 81 adapted to the target distribution.

82 This article is structured as follows. In Section 2 we detail our approach and
 83 theoretical results. In Section 3 we demonstrate how our approach can be implemented
 84 in practice. In Section 4 we give our numerical examples. Section 5 summarizes the
 85 article. The appendix includes the assumptions, technical results and proofs of our
 86 main results.

1.1. Notations. Let $(\mathbf{X}, \mathcal{X})$ be a measurable space. For $\varphi : \mathbf{X} \rightarrow \mathbb{R}$ we write $\mathcal{B}_b(\mathbf{X})$ and $\text{Lip}(\mathbf{X})$ as the collection of bounded measurable and Lipschitz functions respectively. For $\varphi \in \mathcal{B}_b(\mathbf{X})$, we write the supremum norm $\|\varphi\| = \sup_{x \in \mathbf{X}} |\varphi(x)|$. For $\varphi \in \mathcal{B}_b(\mathbf{X})$, $\text{Osc}(\varphi) = \sup_{(x,y) \in \mathbf{X} \times \mathbf{X}} |\varphi(x) - \varphi(y)|$ and we write $\text{Osc}_1(\mathbf{X})$ for the set of functions φ on \mathbf{X} such that $\text{Osc}(\varphi) = 1$. For $\varphi \in \text{Lip}(\mathbf{X})$, we write the Lipschitz constant $\|\varphi\|_{\text{Lip}}$. $\mathcal{P}(\mathbf{X})$ denotes the collection of probability measures on $(\mathbf{X}, \mathcal{X})$. For a measure μ on $(\mathbf{X}, \mathcal{X})$ and a $\varphi \in \mathcal{B}_b(\mathbf{X})$, the notation $\mu(\varphi) = \int_{\mathbf{X}} \varphi(x) \mu(dx)$ is used. Let $K : \mathbf{X} \times \mathcal{X} \rightarrow [0, 1]$ be a Markov kernel and μ be a measure then we use the notations $\mu K(dy) = \int_{\mathbf{X}} \mu(dx) K(x, dy)$ and for $\varphi \in \mathcal{B}_b(\mathbf{X})$, $K(\varphi)(x) = \int_{\mathbf{X}} \varphi(y) K(x, dy)$. For a sequence of Markov kernels K_1, \dots, K_n we write

$$K_{1:n}(x_0, dx_n) = \int_{\mathbf{X}^{n-1}} \prod_{p=1}^n K_p(x_{p-1}, dx_p).$$

87 For $\mu, \nu \in \mathcal{P}(\mathbf{X})$, the total variation distance is written $\|\mu - \nu\|_{\text{tv}} = \sup_{A \in \mathcal{X}} |\mu(A) - \nu(A)|$. For $A \in \mathcal{X}$ the indicator is written $\mathbb{1}_A(x)$. \mathcal{U}_A denotes the uniform distribution
 88 on the set A . $\mathcal{N}(a, b)$ is the one-dimensional Gaussian distribution of mean a and
 89 variance b .

2. Model and Approach. We are given a HMM and we seek to compute

$$\mathbb{E}_{\pi_{n,0}}[\varphi(X_0)|y_{0:n}] = \frac{\int_{\mathbf{X}^{n+1}} \varphi(x_0) \prod_{p=0}^n g(x_p, y_p) f(x_{p-1}, x_p) dx_{0:n}}{\int_{\mathbf{X}^{n+1}} \prod_{p=0}^n g(x_p, y_p) f(x_{p-1}, x_p) dx_{0:n}}$$

where $f(x_{-1}, x_0) := f(x_0)$ and for ease of simplicity we suppose that $\varphi \in \mathcal{B}_b(\mathbf{X}) \cap \text{Lip}(\mathbf{X})$ and \mathbf{X} is a compact subspace of the real line. $\pi_{n,0}$ is the probability density (we also use the same symbol for probability measure) of the smoother given n observations at the co-ordinate at time 0. That is

$$\pi_{n,0}(x_0|y_{0:n}) \propto \int_{\mathbf{X}^n} \prod_{p=0}^n g(x_p, y_p) f(x_{p-1}, x_p) dx_{1:n}.$$

Let $0 < n^* < n$ be fixed, then we propose to consider

$$\mathbb{E}_{\pi_{n^*,0}}[\varphi(X_0)|y_{0:n^*}] = \mathbb{E}_{\pi_{0,0}}[\varphi(X_0)|y_0] + \sum_{p=1}^{n^*} \{\mathbb{E}_{\pi_{p,0}}[\varphi(X_0)|y_{0:p}] - \mathbb{E}_{\pi_{p-1,0}}[\varphi(X_0)|y_{0:p-1}]\}.$$

2.1. Case $X \subset \mathbb{R}$. Let us denote the CDF of $\pi_{p,0}$ as $\Pi_{p,0}$. An approximation of $\mathbb{E}_{\pi_{p,0}}[\varphi(X_0)|y_{0:p}] - \mathbb{E}_{\pi_{p-1,0}}[\varphi(X_0)|y_{0:p-1}]$ is

$$\frac{1}{N_p} \sum_{i=1}^{N_p} [\varphi(\Pi_{p,0}^{-1}(U^i)) - \varphi(\Pi_{p-1,0}^{-1}(U^i))]$$

91 where for $i \in \{1, \dots, N_p\}$, $U^i \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}_{[0,1]}$ and $\Pi_{p,0}^{-1}$ is the (generalized) inverse CDF
 92 of $\Pi_{p,0}$. If we do this independently for each $p \in \{1, \dots, n\}$ and use an independent
 93 estimator $\frac{1}{N_0} \sum_{i=1}^{N_0} \varphi(\Pi_0^{-1}(U^i))$ for $\mathbb{E}_{\pi_{0,0}}[\varphi(X_0)|y_0]$ one can estimate $\mathbb{E}[\varphi(X_0)|y_{0:n}]$.
 94 The utility of the coupling is that it is optimal in terms of 2-Wasserstein distance.
 95 We have the following result, where the assumption and proof are in the appendix.

THEOREM 2.1. *Assume (A1). Then there exists $\rho \in (0, 1)$, $C < +\infty$ such that for any $\varphi \in \mathcal{B}_b(X) \cap \text{Lip}(X)$, $n^* \geq p \geq 1$, $N_p \geq 1$, we have*

$$\text{Var} \left[\frac{1}{N_p} \sum_{i=1}^{N_p} [\varphi(\Pi_{p,0}^{-1}(U^i)) - \varphi(\Pi_{p-1,0}^{-1}(U^i))] \right] \leq \frac{C \rho^{p-1} \|\varphi\|_{\text{Lip}}^2}{N_p}.$$

The main implication of the result is the following. In the approach to be considered later in this paper the cost of computing (an approximation of) $(\Pi_{p,0}^{-1}, \Pi_{p-1,0}^{-1})$ is $\mathcal{O}(1)$ per time step. So the cost of this method is $C(n^* + \sum_{p=0}^{n^*} N_p)$. Thus the MSE and cost associated to this algorithm are (at most in the first case)

$$C(\|\varphi\|^2 \vee \|\varphi\|_{\text{Lip}}^2) \left(\frac{1}{N_0} + \sum_{p=1}^{n^*} \frac{\rho^{p-1}}{N_p} + \rho^{2n} \right)$$

96 and

$$97 \quad (2.1) \quad C(n^* + \sum_{p=0}^{n^*} N_p).$$

Let $\epsilon > 0$ be given. To achieve an MSE of $\mathcal{O}(\epsilon^2)$ we can choose $n^* = \lceil \log(\epsilon) / \log(\rho) \rceil$ (here we of course mean $n^* = \lceil \log(\epsilon) / \log(\rho) \rceil$, but this is omitted for simplicity) and $N_p = \epsilon^{-2}(p+1)^{-1-\delta}$ for any $\delta > 0$ yields that the associated cost is $\mathcal{O}(\epsilon^{-2})$. If one just approximates $\mathbb{E}_{\pi_{n,0}}[\varphi(X_0)|y_{0:n}]$ using

$$\frac{1}{N} \sum_{i=1}^N \varphi(\Pi_{n,0}^{-1}(U^i))$$

then, to achieve an MSE of $\mathcal{O}(\epsilon^2)$ the cost would be $\mathcal{O}(n\epsilon^{-2})$ which is considerably larger if n is large. That is, the cost of the ML approach is essentially $\mathcal{O}(1)$ w.r.t. n . If one stops at $n^* = \lceil \log(\epsilon) / \log(\rho) \rceil$ and uses the estimate

$$\frac{1}{N} \sum_{i=1}^N \varphi(\Pi_{n^*,0}^{-1}(U^i))$$

98 to achieve an MSE of $\mathcal{O}(\epsilon^2)$, the cost is $\mathcal{O}(\epsilon^{-2} |\log(\epsilon)|)$. A similar approach can show
 99 that these results are even true when smoothing for $\mathbb{E}[\varphi(X_{0:k})|y_{0:n}]$ for k fixed (and

100 hence $\mathbb{E}[\varphi(X_{s:s+k})|y_{0:n}]$. The strategy of choosing n^* and $N_{0:n^*}$ detailed above, is
 101 the one used throughout the paper. Note that in practice, we do not know ρ , so we
 102 choose a value such as $\rho = 0.8$ which should lead to an n^* which is large enough. This
 103 is also the reason for setting $N_p = \epsilon^{-2}(p+1)^{-1-\delta}$ and not $N_p = \epsilon^{-2}(\rho^{1/2})^{p-1}$ say.

It is remarked that the compactness of \mathbf{X} could be removed by using Kellerer's extension of the Kantorovich-Rubenstein theorem (see [9] for a summary) and then, given that the latter theory is applicable, to show that there exists a $C < +\infty$, $\rho \in (0, 1)$ such that for any $n^* \geq p \geq 1$

$$\sup_{\varphi \in \text{Lip}_1(\mathbf{X})'} |\mathbb{E}_{\pi_{p,0}}[\varphi(X_0)|y_{0:p}] - \mathbb{E}_{\pi_{p-1,0}}[\varphi(X_0)|y_{0:p-1}]| \leq C\rho^{p-1}$$

104 where $\text{Lip}_1(\mathbf{X})'$ is the collection of functions $\varphi : \mathbf{X} \rightarrow \mathbb{R}$ such that for every $(x, y) \in \mathbf{X}^2$,
 105 $|\varphi(x) - \varphi(y)| \leq |x - y|^2$. This can be achieved using the techniques in [17]. Such
 106 an extension is mainly of a technical nature and is not required in the continuing
 107 exposition. We now establish that the construction here can be extended to the case
 108 $\mathbf{X} \subset \mathbb{R}^d$.

2.2. Case $\mathbf{X} \subset \mathbb{R}^d$. We consider the Knothe-Rosenblatt rearrangement, which is assumed to exist (see e.g. [26]). For simplicity of notation, we set $\mathbf{X} = \mathbf{E}^d$ for some compact $\mathbf{E} \subset \mathbb{R}$. Denote by $\Pi_{p,0}(\cdot|x_{1:j})$ the conditional CDF of $\pi_{p,0}(x_{j+1}|x_{1:j})$ with $1 \leq j \leq d-1$. Note that here we are dealing with the d -dimensional co-ordinate at time zero and we are considering conditioning on the first j of these dimensions. Then to approximate $\mathbb{E}_{\pi_{p,0}}[\varphi(X_0)|y_{0:p}] - \mathbb{E}_{\pi_{p-1,0}}[\varphi(X_0)|y_{0:p-1}]$, sample $U_{1:d}^1, \dots, U_{1:d}^{N_p}$, where for $i \in \{1, \dots, N_p\}$, $U_{1:d}^i \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}_{[0,1]^d}$. Then we have the estimate for $\varphi \in \mathcal{B}_b(\mathbf{X}) \cap \text{Lip}(\mathbf{X})$

$$\frac{1}{N_p} \sum_{i=1}^{N_p} [\varphi(\xi_{p,d}^i) - \varphi(\xi_{p-1,d}^i)]$$

109 where for ease of notation, we have set $\xi_{p,1}^i = \Pi_{p,0}^{-1}(U_1^i)$, (resp. $\xi_{p-1,1}^i = \Pi_{p-1,0}^{-1}(U_1^i)$)
 110 and $\xi_{p,j}^i = (\xi_{p,1}^i, \dots, \xi_{p,j-1}^i, \Pi_{p,0}^{-1}(U_j^i|\xi_{p,j-1}^i))$, $2 \leq j \leq d$, (resp. $\xi_{p-1,j}^i =$
 111 $(\xi_{p-1,1}^i, \dots, \xi_{p-1,j-1}^i, \Pi_{p-1,0}^{-1}(U_j^i|\xi_{p-1,j-1}^i))$, $2 \leq j \leq d$). We have the following result,
 112 whose proof and assumptions are in the appendix.

THEOREM 2.2. *Assume (A1-2). Then there exists $\rho \in (0, 1)$, $C < +\infty$ such that for any $\varphi \in \mathcal{B}_b(\mathbf{X}) \cap \text{Lip}(\mathbf{X})$, $n^* \geq p \geq 1$, $N_p \geq 1$, we have*

$$\mathbb{V}\text{ar} \left[\frac{1}{N_p} \sum_{i=1}^{N_p} [\varphi(\xi_{p,d}^i) - \varphi(\xi_{p-1,d}^i)] \right] \leq \frac{C\rho^{p-1} \|\varphi\|_{\text{Lip}}^2}{N_p}.$$

113 As will be detailed in the following section and in particular in Algorithm 3.1,
 114 it is often more convenient in practice to use the standard normal distribution
 115 instead of the uniform distribution as a base distribution. The only difference is
 116 that samples from the standard normal distribution first have to be mapped through
 117 the corresponding CDF before taking the inverse image through the CDF of interest,
 118 e.g. $\Pi_{p,0}^{-1}(\cdot|x_{1:j})$ for some $p \geq 0$ and some $1 \leq j \leq d-1$.

119 We end this section with some remarks. Firstly, the MLMC strategy could be
 120 debiased w.r.t. the time parameter using the trick in [25], which is a straightforward
 121 extension. One minor issue with this methodology, is that the variance can blow up in
 122 some scenarios. Secondly, the idea of using the approach in [25], when approximating

123 $\mathbb{E}[\varphi(X_{0:n})|y_{0:n}]$ has been adopted in [16]. The authors use a conditional version of the
 124 coupled particle filter (e.g. [5, 18]) to couple smoothers, versus the optimal Wasserstein
 125 coupling. The goal in [16] is unbiased estimation which is complementary to ideas in
 126 this article, where we focus upon reducing the cost of large lag smoothing.

127 3. Transport methodology.

3.1. Standard Approach. The basic principle of the transport methodology introduced in [26] is to determine a mapping T relating a base distribution η , e.g. the normal distribution, to a potentially sophisticated target distribution $\tilde{\pi}$ related to the problem of interest. The distribution η should be easy to sample from so that, given the map T , we can obtain samples from $\tilde{\pi}$ by simply mapping samples from η via T . More precisely, the considered mapping T is characterised by

$$T_{\#}\eta(x) = \eta(T^{-1}(x))|\det \nabla T^{-1}(x)| = \tilde{\pi}(x),$$

that is, the *push-forward* distribution of η by T is $\tilde{\pi}$. Such a mapping can be approximated using deterministic or stochastic optimisation methods. However, the underlying optimisation problem is only amenable when the space on which $\tilde{\pi}$ is defined is of a low dimension, e.g. up to 4. This is not the case in general for the smoothing distributions introduced in the previous sections, especially as the number of observations increases. This is addressed in [26] by identifying the dependence structure between the random variables of interest. In particular, for a hidden Markov model on \mathbb{R}^d , it is possible to decompose the problem into transport maps of dimension $2d$, which does not depend on the number n of observations that define the smoother. The problem at time p can be solved by introducing a mapping T_p of the form

$$T_p(x_p, x_{p+1}) = \begin{bmatrix} T_p^0(x_p, x_{p+1}) \\ T_p^1(x_{p+1}) \end{bmatrix}$$

which will transform the $2d$ -dimensional base distribution η_{2d} into a target distribution related to the considered hidden Markov model, as detailed below. This target distribution can be expressed as

$$\tilde{\pi}_p(x_p, x_{p+1}) \propto \eta_d(x_p) f(T_{p-1}^1(x_p), x_{p+1}) g(x_{p+1}, y_{p+1}),$$

for any $p > 0$, which can be seen to be related to the 1-lag smoother. When $p = 0$, we simply define $\tilde{\pi}_0(x_0, x_1) = f(x_0) f(x_0, x_1) g(x_0, y_0) g(x_1, y_1)$. The base distribution η_{2d} (resp. η_d) is the standard normal distribution of dimension $2d$ (resp. d). The mapping T_p can be embedded into the $2d(n+1)$ -dimensional identity mapping as

$$\bar{T}_p(x_0, \dots, x_n) = (x_0, \dots, x_{p-1}, T_p^0(x_p, x_{p+1}), T_p^1(x_{p+1}), x_{p+2}, \dots, x_n)^t,$$

with \cdot^t denoting the matrix transposition. It follows that

$$\mathbf{T}_n = \bar{T}_0 \circ \dots \circ \bar{T}_n$$

128 is the map such that the pushforward $(\mathbf{T}_n)_{\#}\eta_{d(n+1)}$ is equal to the probability
 129 density function of the smoother at time n . Obtaining samples from the smoothing
 130 distribution is then straightforward: it suffices to sample from $\eta_{d(n+1)}$ and to map the
 131 obtained sample via \mathbf{T}_n .

132 Even in low dimension, the optimisation problem underlying the computation of
 133 the transport maps of interest is not trivial. One first has to consider an appropriate

134 parametrisation of these maps, e.g. via polynomial representations. The parameters
 135 of the considered representation then have to be determined using the following
 136 optimisation problem

$$137 \quad (3.1) \quad T_p^* = \operatorname{argmin}_T \left[\log \tilde{\pi}_p(T(X)) + \log (\det \nabla T(X)) - \log \eta_{2d}(X) \right],$$

where the minimum is taken over the set of monotone increasing lower-triangular maps. This minimisation problem can be solved numerically by considering a parametrised family of maps and deterministic or stochastic optimisation methods. Let T be any acceptable map in the minimisation (3.1) and denote by $T^{(i)}$ the i^{th} component of T , which only depends on the i^{th} first variables, $i \in \{1, \dots, 2d\}$, then the considered parametrisation can be expressed as

$$T^{(i)}(x_1, \dots, x_i) = a_i(x_1, \dots, x_{i-1}) + \int_0^{x_i} b_i(x_1, \dots, x_{i-1}, t)^2 dt$$

138 for some real-valued functions a_i and b_i on \mathbb{R}^{i-1} and \mathbb{R}^i respectively. It is assumed
 139 that the functions $x_j \mapsto a_i(x_1, \dots, x_{i-1})$ and $x_j \mapsto b_i(x_1, \dots, x_{i-1}, t)$ are **probabilists'**
 140 **Hermite** functions [2] extended with constant and linear components for any $j \leq i-1$,
 141 and the function $t \mapsto b_i(x_1, \dots, x_{i-1}, t)$ is also a **probabilists' Hermite** function which
 142 is only extended with a constant component. In particular, these functions take the
 143 form

$$144 \quad a_i(x_1, \dots, x_{i-1}) = \sum_{k=1}^{2d(o_{\text{map}}+1)} c_k \Phi_k(x_1, \dots, x_{i-1})$$

$$145 \quad b_i(x_1, \dots, x_{i-1}, t) = \sum_{k=1}^{2do_{\text{map}}} c'_k \Psi_k(x_1, \dots, x_{i-1}, t)$$

147 with o_{map} the map order, with $\{c_k\}_{k \geq 1}$ and $\{c'_k\}_{k \geq 1}$ some collections of real coefficients
 148 and with Φ_k and Ψ_k basis functions based on the above mentioned **probabilists'**
 149 **Hermite** functions. The expectation in (3.1) is then approximated using a Gauss
 150 quadrature of order o_{exp} in each dimension and the minimisation is solved via the
 151 Newton algorithm using the conjugate-gradient method for each step.

152 The desired function T_p can be recovered through the relation

$$153 \quad (3.2) \quad T_p((x_{p,1}, \dots, x_{p,d}), (x_{p+1,1}, \dots, x_{p+1,d})) =$$

$$154 \quad (S_\sigma \circ T_p^* \circ S_\sigma)(x_{p,1}, \dots, x_{p,d}, x_{p+1,1}, \dots, x_{p+1,d}),$$

157 where $\sigma = (2d, 2d-1, \dots, 1)$ and S_σ is the linear map corresponding to the
 158 permutation matrix of σ , which verifies $S_\sigma^{-1} = S_\sigma$.

159 **3.2. Fixed-Point Smoothing with Transport Maps.** The approach
 160 described in Section 3.1 allows for obtaining samples from the distribution $\pi_{n,0}$ of
 161 X_0 given $(Y_0, \dots, Y_n) = (y_0, \dots, y_n)$ by simply retaining the first d components of
 162 samples from $\eta_{d(n+1)}$ after mapping them through T_n . However, the computational
 163 cost associated with the mapping of samples by T_n increases with n , making the
 164 complexity of the method of the order $\mathcal{O}(n^2)$.

This can however be addressed by considering X_0 as a parameter and by only propagating the transport map corresponding to the posterior distribution of

(X_0, X_n) . This approach has been suggested in [26, section 7.4]. We assume in the remainder of this section that observations start at time step 1 instead of 0. When considering X_0 as a parameter, the elementary transport maps take the form

$$T_p(x_0, x_p, x_{p+1}) = \begin{bmatrix} T_p^{X_0}(x_0) \\ T_p^0(x_0, x_p, x_{p+1}) \\ T_p^1(x_0, x_{p+1}) \end{bmatrix}.$$

and the corresponding target distributions become

$$\tilde{\pi}_1(x_0, x_1, x_2) \propto p_0(x_0)f(x_0, x_1)f(x_1, x_2)g(x_1, y_1)g(x_2, y_2),$$

and

$$\tilde{\pi}_p(x_0, x_p, x_{p+1}) \propto \eta_{2d}(x_0, x_p)f(T_{p-1}^1(x_0, x_p), x_{p+1})g(x_{p+1}, y_{p+1}),$$

for any $p > 1$. The transport map associated with the posterior distribution of (X_0, X_n) is

$$\hat{T}_n(x_0, x_n) = \begin{bmatrix} T_1^{X_0} \circ \dots \circ T_{n-1}^{X_0}(x_0) \\ T_{n-1}^1(x_0, x_n) \end{bmatrix}.$$

165 By recursively approximating the composition $T_1^{X_0} \circ \dots \circ T_{n-1}^{X_0}$ by a single map, the
 166 computation of samples from the posterior distribution of X_0 becomes linear in time.
 167 The pseudo-code for this approach is given in Algorithm 3.1.

168 4. Case Studies.

169 4.1. Linear Gaussian.

170 **4.1.1. Theoretical Result.** The results in Section 2 do not apply to the linear
 171 Gaussian case. We extend our results to this scenario. We assume that the dynamical
 172 and observations models are one-dimensional as well as linear and Gaussian such that
 173 the state and observation random variables at time n can be defined as

$$174 (4.1a) \quad X_n | x_{n-1} \sim \mathcal{N}(\alpha x_{n-1}, \beta^2), \quad n \geq 1$$

$$175 (4.1b) \quad Y_n | x_n \sim \mathcal{N}(x_n, \tau^2), \quad n \geq 0$$

177 and $X_0 \sim \mathcal{N}(0, \sigma^2)$, for some $\alpha \in \mathbb{R}$ and some $\beta, \sigma, \tau > 0$. We have the following
 178 result, whose proof is in the appendix.

THEOREM 4.1. *Assuming that $\text{Var}(X_p | y_{0:p}) \approx \gamma^2$ for all p large enough, it holds that*

$$\text{Var} \left[\frac{1}{N_p} \sum_{i=1}^{N_p} [\Pi_{p,0}^{-1}(U^i) - \Pi_{p-1,0}^{-1}(U^i)] \right] = \mathcal{O} \left(\frac{1}{N_p} \left(\alpha + \frac{\beta^2}{\alpha\gamma^2} \right)^{-2p} \right).$$

179 Theorem 4.1 shows that, under assumptions on the parameters of the model, the
 180 variance of the approximated multilevel term at level p tends to 0 exponentially fast
 181 in p and with an order of $1/N_p$ for the number of samples. This theorem also indicates
 182 that the behaviour depends on all the parameters in the model, although implicitly
 183 in τ . For instance, if $\beta \gg \tau$ then one can consider $\gamma = \tau$ in the above expression. The
 184 assumption about the variance of the filter can be justified in terms of reachability
 185 and observability of the system [20].

186 This rate can get extremely beneficial for the proposed approach when β is large
 187 and γ is small, however it can also make it of little use in the opposite case. This
 188 does not come as a surprise since a large β means that the initial condition is quickly
 189 forgotten so that obtaining a high number of samples from the smoother $\pi_{p,0}$ for large
 190 p would be inefficient, whereas small values of β incur a much higher dependency
 191 between the initial state and the observations at different time steps.

Algorithm 3.1 Multilevel transport

```

1: input:  $\epsilon, \delta, \rho$ 
2: Output: estimate  $\hat{X}_0$  of  $\varphi(X_0) \mid y_{0:n^*}$ 
3:  $n^* = \log(\epsilon) / \log(\rho)$ 
4: for  $p = 1, \dots, n^*$  do
5:   if  $p = 1$  then
6:      $\tilde{\pi}_p(x_0, x_1, x_2) \propto p_0(x_0)f(x_0, x_1)f(x_1, x_2)g(x_1, y_1)g(x_2, y_2)$ 
7:   else
8:      $\tilde{\pi}_p(x_0, x_p, x_{p+1}) \propto \eta_{2d}(x_0, x_p)f(T_{p-1}^1(x_0, x_p), x_{p+1})g(x_{p+1}, y_{p+1})$ 
9:      $\triangleright T_{p-1}^1$  is the second component of  $\hat{T}_{p-1}$ 
10:   end if
11:    $\eta = \mathcal{N}(\mathbf{0}_{2d}, \mathbf{I}_{2d})$ 
12:    $\hat{T}_p = \text{FilteringDistributionTransportMap}(\eta, \tilde{\pi}_p)$ 
13:    $\triangleright$  Compute transport map from  $\eta$  to the law of  $(X_0, X_p) \mid y_{1:p}$  based on  $\tilde{\pi}_p$ 
14:    $N_p = \epsilon^{-2}(p+1)^{-1-\delta}$   $\triangleright$  Compute the number of samples
15:   for  $i = 1, \dots, N_p$  do
16:      $S \sim \eta$ 
17:      $\xi_p^i = \hat{T}_p(S)$ 
18:     if  $p = 1$  then
19:        $\zeta_p^i = \varphi(\xi_p^{i,1:d})$   $\triangleright$  Map the first  $d$  components of  $\xi_p^i$  through  $\varphi$ 
20:     else
21:        $\xi_{p-1}^i = \hat{T}_{p-1}(S)$ 
22:        $\zeta_p^i = \varphi(\xi_p^{i,1:d}) - \varphi(\xi_{p-1}^{i,1:d})$ 
23:     end if
24:   end for
25:    $\hat{X}_0 \leftarrow \hat{X}_0 + \frac{1}{N_p} \sum_{i=1}^{N_p} \zeta_p^i$ 
26: end for

```

4.1.2. Numerical Results. The performance of the proposed method is first assessed in the linear-Gaussian case where an analytical solution of the fixed-point smoothing problem is available, this solution being known as the Rauch-Tung-Striebel smoother [24]. More specifically, we consider the model (4.1) with $X_0 \sim \mathcal{N}(1, \sigma^2)$, $\sigma = 2$ and $\alpha = \beta = \tau = 1$. The transport maps of interest are approximated¹ to the order $o_{\text{map}} = 3$ while the expectation is approximated to the order $o_{\text{exp}} = 5$ and the minimisation is performed with a tolerance of 10^{-4} . The number of samples at each time step as well as the time horizon n^* is computed according to the method proposed in Section 2.1 with different values for the parameter ϵ and with $\rho = 0.8$. The performance of the proposed method is compared against the PaRIS algorithm introduced in [22] using the observations y_1, \dots, y_{50} with a varying number N of samples and with $\tilde{N} = 2$ terms for the propagation of the estimate of X_0 . In the simulations, it always holds that $n^* \leq 50$ to ensure the fairness of the comparison. The criteria for performance assessment is the MSE at the final time step, defined as

$$\frac{1}{M} \sum_{i=1}^M (\hat{x}_i - x^*)^2$$

¹The solver used for the determination of the transport maps is the one provided at <http://transportmaps.mit.edu/docs/index.html>

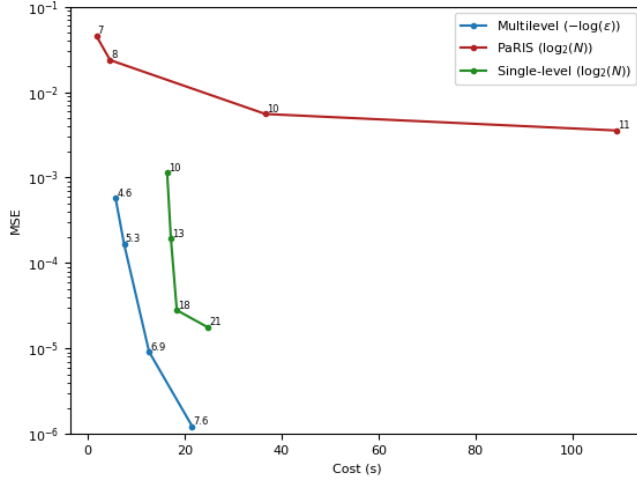


FIG. 1. Performance of the proposed method against the PaRIS algorithm and the single-level transport-map approach for the linear-Gaussian model, averaged over 100 Monte Carlo simulations. The reference for the computation of the MSE is the Rauch-Tung-Striebel smoother. The displayed cost for the multilevel approach includes the computation of the transport maps.

192 where M is the number of Monte Carlo simulations, \hat{x}_i is the estimate of $X_0 \mid y_{1:n^*}$
 193 (with $n^* = 50$ for the PaRIS algorithm) and where x^* is the corresponding estimate
 194 given by the Rauch-Tung-Striebel smoother.

195 The values of the MSE at the final time obtained in simulations are shown
 196 in Figure 1 where the proposed approach displays smaller errors than the PaRIS
 197 algorithm for different values of ϵ and N . The comparison is also made with a single-
 198 level transport-map approach, i.e. without the multilevel decomposition, for different
 199 numbers of samples. The advantage when representing the probability distributions
 200 of interest with transport maps is that the computational effort required to obtain a
 201 sample is extremely limited once the maps have been determined. For instance, the
 202 highest and lowest considered values of ϵ in Figure 1 correspond to $N_1 = 1250$ and
 203 $N_1 = 500,000$ samples respectively, which induces a comparatively small increase in
 204 computational time.

205 In this linear-Gaussian case, using maps of order $o_{\text{map}} < 3$ would have been
 206 sufficient, however this would have been equivalent to making an assumption on
 207 the type of distribution considered for the proposed algorithm whereas the PaRIS
 208 algorithm makes no such assumption. The reason for choosing specifically $o_{\text{map}} = 3$ is
 209 that this value was found to be sufficient for nonlinear models as in the next section.

210 **4.2. Stochastic Volatility Model.** In order to further demonstrate the
 211 performance of the proposed approach, the assessment conducted in the previous
 212 section is applied to the estimation of $X_0 \mid y_{1:n^*}$ in a non-linear case. A stochastic
 213 volatility model is considered with

$$214 \quad X_n = \mu + \phi(X_{n-1} - \mu) + V_n, \quad n \geq 1, \quad X_0 \sim \mathcal{N}\left(\mu, \frac{1}{1 - \phi^2}\right)$$

$$215 \quad Y_n = W_n \exp\left(\frac{1}{2}X_n\right), \quad n \geq 0$$

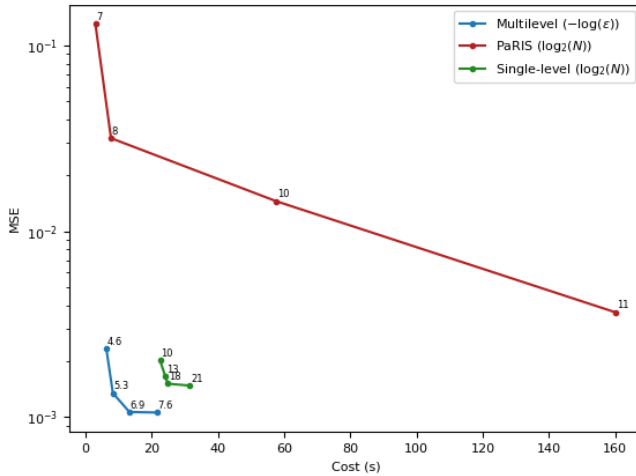


FIG. 2. Performance of the proposed method against the PaRIS algorithm and the single-level transport-map approach for the stochastic volatility model, averaged over 100 Monte Carlo simulations. The reference for the computation of the MSE is the PaRIS algorithm with 2^{14} samples. The displayed cost for the multilevel approach includes the computation of the transport maps.

216 with $V_n \sim \mathcal{N}(0, \beta^2)$ and $W_n \sim \mathcal{N}(0, 1)$, where $\mu = -0.5$, $\phi = 0.95$ and $\beta = 0.25$. In the
 217 absence of an analytical solution, the reference is determined by the PaRIS algorithm
 218 with $N = 2^{14}$ samples. Since the observation process of this model is generally less
 219 informative than the one of the Gaussian model, the PaRIS algorithm is given the
 220 observations up to the time step 50 and, similarly, it is ensured that $n^* \leq 50$ for
 221 the proposed approach. The other parameters are the same as in the linear-Gaussian
 222 case, that is maps of order $o_{\text{map}} = 3$ are used, the expectation is approximated to the
 223 order $o_{\text{exp}} = 5$ and the minimisation is performed with a tolerance of 10^{-4} .

224 The MSE at the final time obtained for the two considered methods is shown
 225 in Figure 2. Once again, the error for the proposed approach is lower than for the
 226 PaRIS algorithm although the difference is less significant. In particular, the gain
 227 in accuracy between the lowest and the second lowest value of ϵ seem to indicate
 228 that simply increasing the number of samples would not allow for reducing the error
 229 much further. However, increasing the order of the transport maps or decreasing
 230 the tolerance in the optimisation could further reduce the error, although with a
 231 significantly higher computational cost.

232 The computational costs obtained for the two models considered in simulations
 233 are shown in Figure 3 for different values of ϵ . These results confirm the order $\mathcal{O}(\epsilon^{-2})$
 234 that was predicted in Section 2.

235 **5. Summary.** In this article we have considered large lag smoothing for HMMs,
 236 using the MLMC method. We showed that under an optimal coupling when the
 237 hidden state is in dimension 1 or higher, but on a compact space that, essentially,
 238 the cost can be decoupled from the time parameter of the smoother. As this optimal
 239 method is not possible in practice, we showed how it could be approximated and
 240 established numerically that our theory still holds in this approximated case. Several
 241 extensions to the work are possible. Firstly, to extend our theoretical results to the

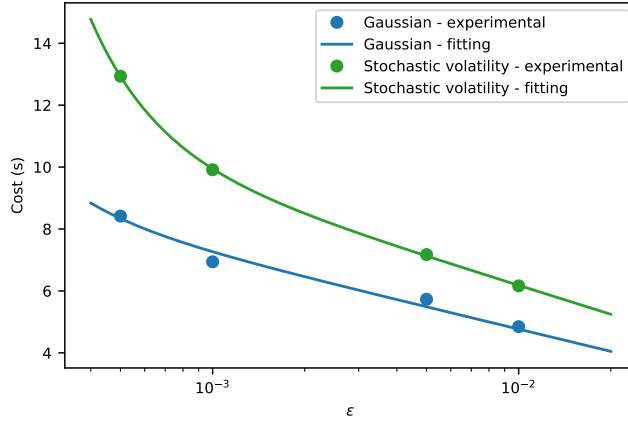


FIG. 3. Computational cost as a function of ϵ , averaged over 100 Monte Carlo simulations. The fitted curves are based on a function of the form $\epsilon \mapsto -a\epsilon^{-2} - b\log(\epsilon)$, with a and b some parameters, which is justified by the form of the cost (2.1).

242 case of the approximated coupling. Secondly, to investigate whether the coupling used
 243 in [16] can also yield, theoretically, the same improvements that have been seen in the
 244 work in this article.

245 **Appendix A. Variance Proofs.** We write the density (or probability measure)
 246 of the smoother, at time p , on the co-ordinate at time zero as $\pi_{p,0}$ and the associated
 247 CDF as $\Pi_{p,0}$ (with generalized inverse $\Pi_{p,0}^{-1}$). Recall that throughout \mathbf{X} is a compact
 248 subspace of \mathbb{R}^d . Throughout the observations are fixed and often omitted from the
 249 notations. The appendix gives our main assumptions, followed by a technical Lemma
 250 (Lemma A.1) which features some technical results used in the proofs. Then the proof
 251 of Theorem 2.1 is given. The appendix is concluded by a second technical Lemma
 252 (Lemma A.2) followed by the proof of Theorem 2.2.

253 **(A1)** There exists $0 < \underline{C} < \overline{C} < +\infty$ such that

$$254 \quad \inf_{x \in \mathbf{X}} g(x, y_0) f(x) \wedge \inf_{p \geq 1} \inf_{(x, x') \in \mathbf{X}^2} g(x', y_p) f(x, x') \geq \underline{C}$$

$$255 \quad \sup_{x \in \mathbf{X}} g(x, y_0) f(x) \vee \sup_{p \geq 1} \sup_{(x, x') \in \mathbf{X}^2} g(x', y_p) f(x, x') \leq \overline{C}.$$

256 **(A2)** There exists $C < +\infty$ such that for every $(x, x') \in \mathbf{X}^2$

$$257 \quad |g(x, y_0) - g(x', y_0)| \leq C|x - x'|$$

$$258 \quad \sup_{z \in \mathbf{X}} |f(x, z) - f(x', z)| \leq C|x - x'|$$

$$259 \quad |f(x) - f(x')| \leq C|x - x'|.$$

260 Below $\pi_{p,0}(\cdot | x_{1:j})$ denotes the probability of the $(j+1)^{th}$ co-ordinate of the
 261 smoother at time 0, given the first j -co-ordinates at time 0, and conditional upon
 262 the observations up-to time p .

263 **LEMMA A.1.** Assume (A1-2). Then there exists $(C, C') \in (0, \infty)^2$, $\rho \in (0, 1)$ such
 264 that

- 265 1. for any $1 \leq j \leq d$, $\sup_{p \geq 0} \pi_{p,0}(x_{0,1:j}) \leq C$, $\inf_{p \geq 0} \pi_{p,0}(x_{0,1:j}) \geq C'$
 266 2. for any $p \geq 1$, $\|\pi_{p,0} - \pi_{p-1,0}\|_{\text{tv}} \leq C\rho^{p-1}$
 267 3. for any $1 \leq j \leq d$, $p \geq 1$, $\sup_{x_{1:j} \in \mathbf{E}^j} \|\pi_{p,0}(\cdot|x_{1:j}) - \pi_{p-1,0}(\cdot|x_{1:j})\|_{\text{tv}} \leq C\rho^{p-1}$
 268 4. for any $p \geq 0$, $(x, x') \in \mathbf{X}^2$, $|\pi_{p,0}(x) - \pi_{p,0}(x')| \leq C|x - x'|$
 269 5. for any $p \geq 0$, $1 \leq j \leq d$, $(x_{1:j}, x'_{1:j}) \in (\mathbf{E}^j)^2$, $|\pi_{p,0}(x_{1:j}) - \pi_{p,0}(x'_{1:j})| \leq$
 270 $C|x_{1:j} - x'_{1:j}|$.

271 *Proof.* 1. follows trivially from (A1) and the compactness of \mathbf{E} . 2. follows from
 272 the backward Markov chain representation of the smoother and (A1); see for instance
 273 [4] and the references therein.

3. to prove this result, we first consider controlling for any fixed $1 \leq j \leq d$ $p \geq 1$,

$$|\pi_{p,0}(x_{1:j}) - \pi_{p-1,0}(x_{1:j})|.$$

Denoting $\pi_{(p)}$ as the filter at time p and setting for $k \geq 0$

$$B_k(x_{k+1}, x_k) = \frac{\pi_{(k)}(x_k) f(x_k, x_{k+1})}{\int_{\mathbf{X}} \pi_{(k)}(x_k) f(x_k, x_{k+1}) dx_k}$$

274 we can write

275

$$(A.1) \quad |\pi_{p,0}(x_{1:j}) - \pi_{p-1,0}(x_{1:j})| =$$

277

278

$$\text{Osc}(B_0(\cdot, x_{1:j})) \left| [\pi_{(p)} B_{p-1} - \pi_{(p-1)}](B_{p-2:1}) \left(\frac{B_0(\cdot, x_{1:j})}{\text{Osc}(B_0(\cdot, x_{1:j}))} \right) \right|.$$

Using standard results for the total variation distance

$$|\pi_{p,0}(x_{1:j}) - \pi_{p-1,0}(x_{1:j})| \leq \text{Osc}(B_0(\cdot, x_{1:j})) \prod_{s=1}^{p-2} \omega(B_s)$$

279 where $\omega(B_s)$ is the Dobrushin coefficient of the Markov kernel B_s . Standard
 280 calculations yield that there exists a $\rho \in (0, 1)$ such that $\text{Osc}(B_0(\cdot, x_{1:j})) \vee \omega(B_s) \leq C\rho$,
 281 where C does not depend upon $x_{1:j}$. Hence we have shown that

$$(A.2) \quad \sup_{x_{1:j} \in \mathbf{E}^j} |\pi_{p,0}(x_{1:j}) - \pi_{p-1,0}(x_{1:j})| \leq C\rho^{p-1}.$$

283 To prove the result of interest we have for any $\varphi \in \text{Osc}_1(\mathbf{E})$

$$284 \quad |\pi_{p,0}(\varphi|x_{1:j}) - \pi_{p-1,0}(\varphi|x_{1:j})| = \frac{1}{\pi_{p,0}(x_{1:j-1})} \int_{\mathbf{E}} \varphi(x_j) [\pi_{p,0}(x_{1:j}) - \pi_{p-1,0}(x_{1:j})] dx_j +$$

$$285 \quad \frac{\pi_{p-1,0}(x_{1:j-1}) - \pi_{p,0}(x_{1:j-1})}{\pi_{p,0}(x_{1:j-1}) \pi_{p-1,0}(x_{1:j-1})} \int_{\mathbf{E}} \varphi(x_j) \pi_{p-1,0}(x_{1:j}) dx_j.$$

286 The conclusion then follows by using (A.2) and 1..

287 4. follows almost immediately from (A2) and the definition of the smoother. 5.
 288 follows from 4. on marginalization and the compactness of \mathbf{E} . \square

Proof of Theorem 2.1. Standard calculations for i.i.d. random variables and the Lipschitz property of φ clearly yields:

$$\text{Var} \left[\frac{1}{N_p} \sum_{i=1}^N [\varphi(\Pi_{p,0}^{-1}(U^i)) - \varphi(\Pi_{p-1,0}^{-1}(U^i))] \right] \leq \frac{\|\varphi\|_{\text{Lip}}^2}{N_p} \int_{[0,1]} |\Pi_{p,0}^{-1}(u) - \Pi_{p-1,0}^{-1}(u)|^2 du.$$

Now we note that

$$\int_{[0,1]} |\Pi_{p,0}^{-1}(u) - \Pi_{p-1,0}^{-1}(u)|^2 du = W_2(\pi_{p,0}, \pi_{p-1,0})^2$$

where $W_2(\pi_{p,0}, \pi_{p-1,0})$ is the 2-Wasserstein distance between $\pi_{p,0}$ and $\pi_{p-1,0}$. As X is compact it follows

$$W_2(\pi_{p,0}, \pi_{p-1,0})^2 \leq \left(\int_{\mathsf{X}} dx \right)^2 \|\pi_{p,0} - \pi_{p-1,0}\|_{\text{tv}}$$

where $\|\cdot\|_{\text{tv}}$ is the total variation distance. Under our assumptions one can show that there exists $\rho \in (0, 1)$, $C < +\infty$ such that for any $p \geq 1$ (see Lemma A.1 2., which holds when $d = 1$)

$$\|\pi_{p,0} - \pi_{p-1,0}\|_{\text{tv}} \leq C\rho^{p-1}.$$

289 The proof is then easily concluded. \square

LEMMA A.2. *Assume (A1-2). Then there exists $C < +\infty$, $\rho \in (0, 1)$ such that for any $p \geq 1$*

$$\mathbb{E}[|\xi_{p,d}^1 - \xi_{p-1,d}^1|^2] \leq C\rho^{p-1}.$$

290 *Proof.* The proof is by induction on d , the case $d = 1$ being proved by the approach
291 in the proof of Theorem 2.1. Throughout C is a finite constant whose value may
292 change from line-to-line, but does not depend upon p .

293 We suppose the result for $d - 1$ and consider d . For simplicity of notation, we
294 drop the superscript 1 from the notation, e.g. we write $\xi_{p,d}$ instead of $\xi_{p,d}^1$. We have

$$\begin{aligned} 295 \quad \mathbb{E}[|\xi_{p,d} - \xi_{p-1,d}|^2] &= \mathbb{E}[\mathbb{E}[|\xi_{p,d}^1 - \xi_{p-1,d}^1|^2 | U_{1:d-1}]] \\ 296 \quad (\text{A.3}) \quad &\leq C\mathbb{E}[\|\pi_{p,0}(\cdot|\xi_{p,d-1}) - \pi_{p-1,0}(\cdot|\xi_{p-1,d-1})\|_{\text{tv}}] \end{aligned}$$

297 where, to go to the second line, we have used (conditional upon $U_{1:d}$) the relationship
298 between the squared 2-Wasserstein distance and the (generalized) inverse CDF, along
299 with the total variation bound as used in the proof of Theorem 2.1.

Now, we have

$$\begin{aligned} 300 \quad &\|\pi_{p,0}(\cdot|\xi_{p,d-1}) - \pi_{p-1,0}(\cdot|\xi_{p-1,d-1})\|_{\text{tv}} \leq \\ 301 \quad (\text{A.4}) \quad &\|\pi_{p,0}(\cdot|\xi_{p,d-1}) - \pi_{p-1,0}(\cdot|\xi_{p,d-1})\|_{\text{tv}} + \|\pi_{p-1,0}(\cdot|\xi_{p,d-1}) - \pi_{p-1,0}(\cdot|\xi_{p-1,d-1})\|_{\text{tv}}. \end{aligned}$$

302 By Lemma A.1 3. it follows that

$$303 \quad (\text{A.5}) \quad \|\pi_{p,0}(\cdot|\xi_{p,d-1}) - \pi_{p-1,0}(\cdot|\xi_{p,d-1})\|_{\text{tv}} \leq C\rho^{p-1}$$

304 so we consider $\|\pi_{p-1,0}(\cdot|\xi_{p,d-1}) - \pi_{p-1,0}(\cdot|\xi_{p-1,d-1})\|_{\text{tv}}$. For any $\varphi \in \text{Osc}_1(\mathsf{E})$

305

$$306 \quad \pi_{p,0}(\varphi|\xi_{p,d-1}) - \pi_{p-1,0}(\varphi|\xi_{p,d-1}) =$$

$$307 \quad \frac{1}{\pi_{p-1,0}(\xi_{p-1,d-1})} \int_{\mathsf{E}} \varphi(x) [\pi_{p-1,0}(\xi_{p,d-1}, x) - \pi_{p-1,0}(\xi_{p-1,d-1}, x)] dx +$$

$$308 \quad \frac{\pi_{p-1,0}(\xi_{p-1,d-1}) - \pi_{p-1,0}(\xi_{p,d-1})}{\pi_{p-1,0}(\xi_{p,d-1})\pi_{p-1,0}(\xi_{p-1,d-1})} \int_{\mathsf{E}} \varphi(x) \pi_{p-1,0}(\xi_{p-1,d-1}, x) dx.$$

309

310 Applying Lemma A.1 4. to the first term on the R.H.S. and Lemma A.1 5. to the
 311 second term on the R.H.S. along with the boundedness of φ and compactness of \mathbf{E} ,
 312 we have that

$$313 \quad |\pi_{p,0}(\varphi|\xi_{p,d-1}) - \pi_{p-1,0}(\varphi|\xi_{p,d-1})| \leq \frac{C}{\pi_{p-1,0}(\xi_{p-1,d-1})} |\xi_{p,d-1} - \xi_{p-1,d-1}| +$$

$$314 \quad \frac{C}{\pi_{p-1,0}(\xi_{p,d-1})\pi_{p-1,0}(\xi_{p-1,d-1})} |\xi_{p,d-1} - \xi_{p-1,d-1}|.$$

315 Applying Lemma A.1 1. we can then establish that

$$316 \quad (\text{A.6}) \quad \|\pi_{p-1,0}(\cdot|\xi_{p,d-1}) - \pi_{p-1,0}(\cdot|\xi_{p-1,d-1})\|_{\text{tv}} \leq C|\xi_{p,d-1} - \xi_{p-1,d-1}|.$$

Combining (A.5) and (A.6) with (A.4) and noting (A.3), we have shown that

$$\mathbb{E}[|\xi_{p,d} - \xi_{p-1,d}|^2] \leq C\left(\rho^{p-1} + \mathbb{E}[|\xi_{p,d-1} - \xi_{p-1,d-1}|\right).$$

317 The proof is completed by using the Jensen inequality and the induction hypothesis. \square

Proof of Theorem 2.2. We have

$$\text{Var}\left[\frac{1}{N_p} \sum_{i=1}^{N_p} [\varphi(\xi_{p,d}^i) - \varphi(\xi_{p-1,d}^i)]\right] \leq \frac{\|\varphi\|_{\text{Lip}}^2}{N_p} \mathbb{E}[|\xi_{p,d}^1 - \xi_{p-1,d}^1|^2].$$

318 The proof is then completed by applying Lemma A.2. \square

319 **Appendix B. Linear Gaussian Result.**

320 *Proof of Theorem 4.1.* The Rauch-Tung-Striebel smoother gives an expression of
 321 the smoothed mean $m_{p|n}$ and variance $v_{p|n}$ at time p given the observations y_0, \dots, y_n
 322 as

$$323 \quad m_{p|n} = m_{p|p} + c_p(m_{p+1|n} - m_{p+1|p})$$

$$324 \quad v_{p|n} = v_{p|p} + c_p^2(v_{p+1|n} - v_{p+1|p}),$$

with $c_p = \alpha m_{p|p}/m_{p+1|p}$, where $m_{p+1|p}$ and $v_{p+1|p}$ are the predicted mean and
 variance at time $p+1$ given the observations y_0, \dots, y_p . It follows that the mean m_p
 and variance v_p of $\pi_{p,0}$ satisfy similar relations to the filtered means and variances:

$$m_p = \sum_{i=0}^p m_{i|i} \alpha^i (1 - \mathbb{I}_{i < p} \alpha^2 d_p) \prod_{j=0}^{i-1} d_j \quad \text{and} \quad v_p = \sum_{i=0}^p v_{i|i} \alpha^{2i} (1 - \mathbb{I}_{i < p} \alpha^4 d_p^2) \prod_{j=0}^{i-1} d_j^2,$$

where $d_p = v_{p|p}/v_{p+1|p}$ and where \mathbb{I}_c is the indicator of condition c . The objective is
 to compute the order of

$$\Pi_{p,0}^{-1}(u) - \Pi_{p-1,0}^{-1}(u) = m_p - m_{p-1} + \sqrt{2} \text{erf}^{-1}(2u - 1)(\sigma_p - \sigma_{p-1})$$

where $\sigma_p = \sqrt{v_p}$. From the above expression, it follows easily that

$$m_p - m_{p-1} = \alpha^p (m_{p|p} - m_{p|p-1}) \prod_{i=0}^{p-1} d_i \quad \text{and} \quad v_p - v_{p-1} = \alpha^{2p} (v_{p|p} - v_{p|p-1}) \prod_{i=0}^{p-1} d_i^2.$$

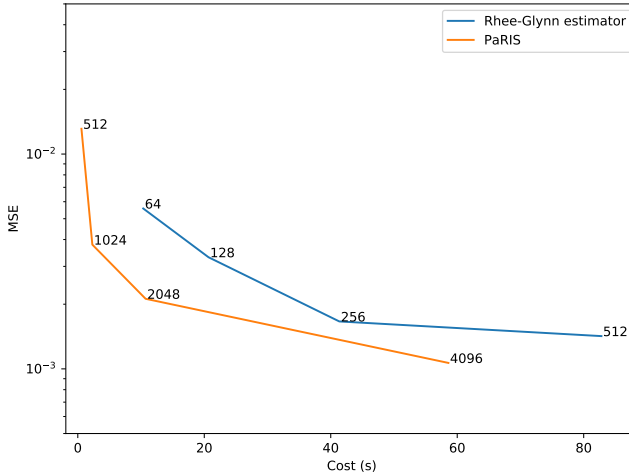


FIG. 4. Performance of the Rhee-Glynn estimator against the PaRIS algorithm with a linear-Gaussian model, averaged over 100 Monte Carlo simulations, where the number of samples is indicated on the figure. The reference for the computation of the MSE is the Rauch-Tung-Striebel smoother. The results for the Rhee-Glynn estimator are averaged over 2^{10} runs of the estimator.

which yields the same order for both $m_p - m_{p-1}$ and $\sigma_p - \sigma_{p-1}$. The desired result follows from the fact that

$$\alpha^p \prod_{i=0}^{p-1} d_i = \alpha^p \prod_{i=0}^{p-1} \frac{v_{i|i}}{\alpha v_{i|i} + \beta^2} = \prod_{i=0}^{p-1} \frac{\alpha}{\alpha^2 + \beta^2 / v_{i|i}} = \prod_{i=0}^{p-1} \left(\alpha + \frac{\beta^2}{\alpha v_{i|i}} \right)^{-1},$$

and from the assumption that $v_{p|p} = \text{Var}(X_p | y_{0:p}) \approx \gamma^2$ for all p large enough. \square

Appendix C. The Rhee-Glynn smoothing estimator. We compare the so-called *Rhee-Glynn smoothing estimator* described in [16] with the PaRIS algorithm [22] on the linear-Gaussian model considered in Section 4.1.2. The Rhee-Glynn smoothing estimator is implemented with ancestor sampling [21] and where all the generated paths are used in the estimate of $X_0 | y_{1:n^*}$, as originally suggested in [1] in the context of particle Markov chain Monte Carlo.

The result of the comparison is given in Figure 4 where it appears that the PaRIS algorithm slightly outperforms the Rhee-Glynn smoothing estimator. Although the scenario considered here is linear and Gaussian, none of the compared methods relies on these assumptions so that the conclusions made for this case are generalisable to some other classes of scenarios. This justifies the sole use of the PaRIS algorithm in Section 4 for comparison against the proposed approach.

REFERENCES

- [1] ANDRIEU, C., DOUCET, A. & HOLENSTEIN, R. (2010). Particle Markov chain Monte Carlo methods (with discussion). *J. R. Statist. Soc. Ser. B*, **72**, 269–342.
- [2] BOYD, J.P. (2001). *Chebyshev and Fourier spectral methods*. Courier Corporation.
- [3] BUNCH, P., & GODSILL, S. (2016). Approximations of the optimal importance density using Gaussian particle flow importance sampling. *Journal of the American Statistical Association*, **111**(514), 748–762.

- 345 [4] CAPPÉ, O., RYDEN, T. & MOULINES, É. (2005). *Inference in Hidden Markov Models*. Springer:
346 New York.
- 347 [5] CHOPIN, N. & SINGH, S. S. (2015). On particle Gibbs sampling. *Bernoulli*, **21**, 1855–1883.
- 348 [6] DAUM, F., & HUANG, J. (2008). Particle flow for nonlinear filters with log-homotopy. In *Signal*
349 *and Data Processing of Small Targets*. International Society for Optics and Photonics.
- 350 [7] DEL MORAL, P., DOUCET, A. & SINGH, S. S. (2010). A backward interpretation of Feynman-
351 Kac formulae. *M2AN*, **44**, 947–975.
- 352 [8] DOUCET, A. & JOHANSEN, A. (2011). A tutorial on particle filtering and smoothing: Fifteen
353 years later. In *Handbook of Nonlinear Filtering* (eds. D. Crisan & B. Rozovsky), Oxford
354 University Press: Oxford.
- 355 [9] EDWARDS, D. A. (2011). On the Kantorovich-Rubinstein theorem. *Expos. Math.*, **29**, 387–398.
- 356 [10] GILES, M. B. (2008). Multilevel Monte Carlo path simulation. *Op. Res.*, **56**, 607–617.
- 357 [11] GILES, M. B. (2013). Multilevel Monte Carlo methods. In *Monte Carlo and Quasi-Monte Carlo*
358 *Methods 2012* (pp. 83-103). Springer, Berlin, Heidelberg.
- 359 [12] GREGORY, A., COTTER, C. J., & REICH, S. (2016). Multilevel ensemble transform particle
360 filtering. *SIAM J. Sci. Comput.*, **38**(3), A1317–A1338.
- 361 [13] HEINRICH, S. (2001). Multilevel Monte Carlo methods. In *Large-Scale Scientific Computing*,
362 (eds. S. Margenov, J. Wasniewski & P. Yalamov), Springer: Berlin.
- 363 [14] HENG, J., DOUCET, A., & POKERN, Y. (2015). Gibbs flow for approximate transport with
364 applications to Bayesian computation. arXiv preprint arXiv:1509.08787.
- 365 [15] HOEL, H., LAW, K. J., & TEMPONE, R. (2016). Multilevel ensemble Kalman filtering. *SIAM J.*
366 *Numer. Anal.*, **54**(3), 1813–1839.
- 367 [16] JACOB, P., LINDSTEN, F. & SCHÖN, T. (2017). Smoothing with Couplings of Conditional Particle
368 Filters. arXiv preprint, arXiv:1701.02002.
- 369 [17] JASRA, A. (2015). On the behaviour of the backward interpretation of Feynman-Kac formulae
370 under verifiable conditions. *J. Appl. Probab.*, **52**, 339–359.
- 371 [18] JASRA, A., KAMATANI, K., LAW, K. J. H. & ZHOU, Y. (2017). Multilevel particle filters. *SIAM*
372 *J. Numer. Anal.*, **55**, 3068–3096.
- 373 [19] KANTAS, N., DOUCET, A., SINGH, S. S., MACIEJOWSKI, J. M. & CHOPIN, N. (2015) On Particle
374 Methods for Parameter Estimation in General State-Space Models. *Statist. Sci.*, **30**, 328–
375 351.
- 376 [20] KUMAR, P.R., & PRAVIN V. (1986). *Stochastic systems: Estimation, identification, and*
377 *adaptive control*. Prentice-Hall.
- 378 [21] LINDSTEN, F., JORDAN, M. I., & SCHÖN, T. B. (2014). Particle Gibbs with ancestor sampling.
379 *The Journal of Machine Learning Research*, **15**(1), 2145–2184.
- 380 [22] OLSSON, J. & WESTERBORN, J. (2017). Efficient particle-based online smoothing in general
381 hidden Markov models: The PaRIS algorithm. *Bernoulli*, **23**, 1951–1996.
- 382 [23] PARNO, M., MOSELHY, T., & MARZOUK, Y. (2016). A multiscale strategy for Bayesian inference
383 using transport maps. *SIAM/ASA Journal on Uncertainty Quantification*, **4**(1), 1160–
384 1190.
- 385 [24] RAUCH, H. E., STRIEBEL, C. & TUNG, F. (1965) Maximum likelihood estimates of linear
386 dynamical systems. *AIAA J.* **3**, 1445–1450.
- 387 [25] RHEE, C. H., & GLYNN, P. W. (2015). Unbiased estimation with square root convergence for
388 SDE models. *Op. Res.*, **63**, 1026–1043.
- 389 [26] SPANTINI, A., BIGONI, D. & MARZOUK, Y. (2018). Inference via low-dimensional couplings. *The*
390 *Journal of Machine Learning Research*, **19**(1), 2639–2709.