# A SPLITTING METHOD TO REDUCE MCMC VARIANCE

BY ROBERT J. WEBBER[1], DAVID ARISTOFF[2] AND GIDEON SIMPSON[3]

[1]*Courant Institute of Mathematical Sciences, rw2515@nyu.edu*

[2]*Colorado State University, aristoff@math.colostate.edu*

[3]*Drexel University, grs53@drexel.edu*

We explore whether splitting and killing methods can improve the accuracy of Markov chain Monte Carlo (MCMC) estimates of rare event probabilities, and we make three contributions. First, we prove that "weighted ensemble" is the only splitting and killing method that provides asymptotically consistent estimates when combined with MCMC. Second, we prove a lower bound on the asymptotic variance of weighted ensemble's estimates. Third, we give a constructive proof and numerical examples to show that weighted ensemble can approach this optimal variance bound, in many cases reducing the variance of MCMC estimates by multiple orders of magnitude.

**1. Introduction.** Markov chain Monte Carlo (MCMC) is a stochastic method that empowers researchers to calculate statistics of high-dimensional systems that could not be calculated by other means. The MCMC approach for calculating an integral $\mu(f) = \int \mu(dx) f(x)$ involves simulating a Markov chain $X_t$ that is ergodic with respect to $\mu$ and then forming the trajectory average

$$(1.1) \qquad \mu(f) \approx \frac{1}{T} \sum_{t=0}^{T-1} f(X_t).$$

Here, we use a broader definition of MCMC than is typical. We refer to MCMC as any scheme that computes ergodic averages by simulating a Markov chain and then taking trajectory averages. Our definition thus includes traditional MCMC samplers, such as the random walk Metropolis [30] or Gibbs sampler [20], that require specifying a target density known up to a normalization constant. Our definition also extends to samplers where the ergodic distribution is unknown and can only be ascertained through simulations (e.g., [24, 12]).

Despite the many benefits of MCMC, the approach often performs poorly when estimating probabilities of rare sets. When calculating a small probability $\mu(A) \ll 1$, MCMC requires a long simulation time to ensure accuracy, and running a simulation for such a long time can be prohibitively computationally expensive. This limitation makes MCMC difficult to apply in impactful rare event estimation problems where accurate computations are greatly needed.

In this work, we investigate the possibility of incorporating *splitting and killing* into MCMC to better compute small probabilities. Splitting and killing (commonly abbreviated "splitting") is an approach in which we simulate a collection of Markov chains ("particles") using a common transition kernel $K$. Periodically, we replicate some of the particles to promote progress toward a rare outcome and randomly kill other particles to prevent a population explosion.

Splitting is one of the most prevalent Monte Carlo approaches for rare event sampling [35]. This approach has been developed over seventy years of applications [33, 22, 21, 16, 8],

originally stemming from an idea proposed by John von Neummann in the 1940s [27]. Given this long history and the method's demonstrated track record of success [9, 25], we sought to apply splitting to improve MCMC's accuracy for rare event probability estimation.

However, we acknowledge two factors separating MCMC from splitting as traditionally applied. First, an MCMC method continues for as long as necessary to ensure robust estimates, whereas a splitting method typically ends as soon as the particles reach a predetermined stopping time [29, 36]. Second, an MCMC method uses every data point to compute time averages (1.1), whereas many splitting methods use only the particles' locations at the final algorithmic step [35] to compute estimates.

In this work, we consider a nontraditional approach to splitting that incorporates an arbitrarily long run time and time-averaged estimates. We ask, what happens if we run an ensemble of ergodic Markov chains and apply splitting at regular intervals? As time goes on, can splitting improve the accuracy of MCMC estimates?

Through numerical analysis and computational examples, we begin to answer these questions. We show as $T \to \infty$ many splitting methods experience a catastrophic shrinking of statistical weights, causing all estimates to converge to zero. Moreover, under mild assumptions, we prove that the *only* splitting method providing consistent estimates as $T \to \infty$ is the weighted ensemble (WE) method proposed by Huber and Kim in 1997 [24].

Unique among splitting methods, WE is characterized by a binning procedure applied at every splitting step. The particles are divided into bins, the population in some of the bins is increased through splitting, and the population in other bins is decreased through killing. During this splitting step, at least one particle must remain in each bin, as shown below in Figure 1 below.
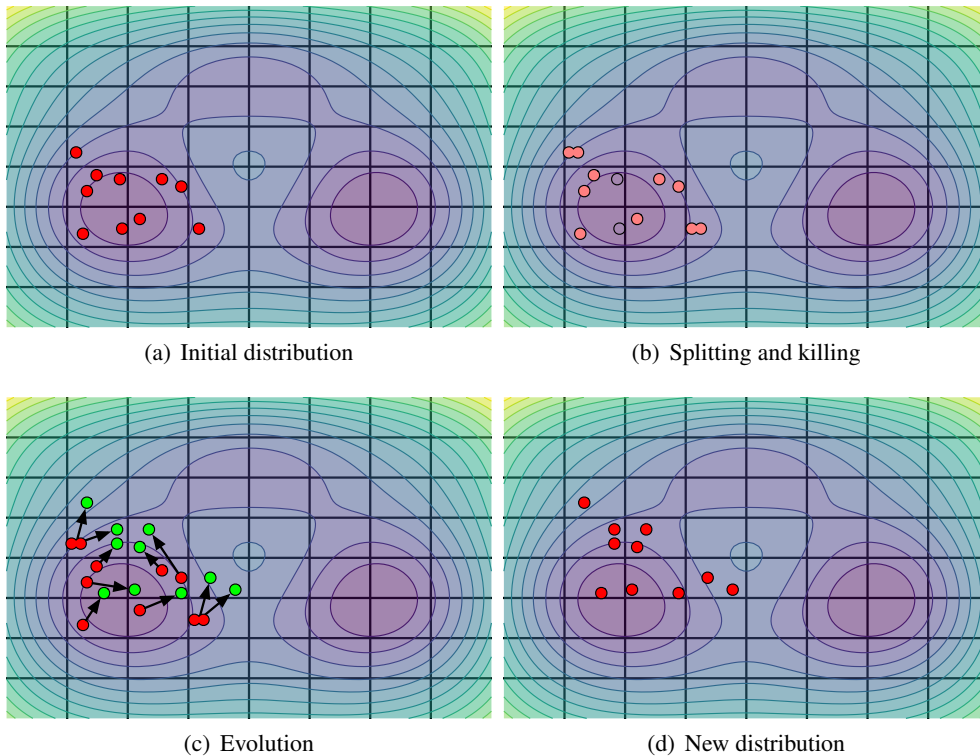


(a) Initial distribution

(b) Splitting and killing

(c) Evolution

(d) New distribution

FIG 1. *WE schematic with Cartesian bins and particles diffusing on an energy landscape.*

Researchers have applied WE with increasing frequency over the past decade [7, 42], and this growing record of applications demonstrates key features that make WE useful:

1. WE relies on forward simulations of a process without introducing bias [41]. This makes WE a non-intrusive method suitable for use with any black-box MCMC sampler. Thus, WE is used in non-equilibrium statistical mechanics, where the invariant distribution is unknown and traditional samplers would not be applicable [24, 12, 42].
2. WE only requires storage at the precise times of splitting/killing. Otherwise, WE avoids recording the complete history of the system, which could be a costly procedure in a high-dimensional state space [17].

As WE has become popular, there have been efforts toward extracting the method's mathematical properties. Zhang and coauthors [41] established the bias properties of WE estimates, Aristoff [2] established the convergence of WE estimates as $T \to \infty$, and Aristoff and Zuckerman [1, 4] developed strategies toward algorithmic optimization. In this past work, however, it was not reported that WE is the *unique* splitting scheme providing consistent estimates as $T \to \infty$. Additionally, major questions remain open about WE's efficiency relative to direct MCMC sampling. In particular, we ask: when does WE produce more accurate estimates than MCMC? Also, what is the lowest possible variance that WE estimates can exhibit?

To help answer these questions, we investigate the asymptotic variance of WE estimates as $T \to \infty$. We establish a lower bound on WE's asymptotic variance that is valid for any number of particles $N$, and we prove that WE can come arbitrarily close to achieving this optimal variance bound as $N \to \infty$. Additionally, we present examples of rare event estimation problems where WE reduces MCMC's asymptotic variance by multiple orders of magnitude. In our numerical experiments, by incorporating sufficiently many particles and optimizing bin allocations, we obtain nearly the optimal variance reduction that WE offers.

Taken as the whole, the impact of our work is both computational and mathematical. On the computational side, we demonstrate how an optimized WE method gives accurate rare event probability estimates that would be extremely costly to obtain by direct MCMC sampling. On the mathematical side, we show that despite the apparent complexity of a splitting scheme's dynamics, the stability and variance properties are governed by simple, fundamental considerations in the limit as $T \to \infty$.

In the rest of this introductory section, we describe our contributions in greater detail and we lay out the plan for the rest of the paper.

1.1. *Ergodicity theory for splitting schemes.* Our first contribution is to explain how a splitting scheme's *statistical weights* influence the method's long-time stability. These statistical weights, which we denote $w_t^1, \ldots, w_t^{N_t}$, are assigned to each of the sampled particles $\xi_t^1, \ldots, \xi_t^{N_t}$. Throughout the scheme, the weights are adjusted in inverse proportion to the amount of splitting and killing that occurs. For example, if a particle is split into two copies, each child receives half the weight of the parent. Conversely, if two equally weighted particles are randomly reduced to a single particle, the surviving particle's weight is doubled. The weights ensure that the splitting scheme can produce estimates

$$(1.2) \qquad \mu(f) \approx \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right),$$

and these estimates are asymptotically unbiased as $T \to \infty$.

While the statistical weights ensure that a splitting method's estimates are asymptotically unbiased, the weights themselves can still degenerate. In experiments, we find that statistical weights converge to zero in many splitting methods, which prevents any possibility of

consistent estimation. Moreover, we give a simple explanation for the shrinking weights by identifying the sum of weights as a nonnegative martingale. A nonnegative martingale must converge to a positive number or to zero in the limit as $T \to \infty$ [28]. Therefore, when the sum of the weights fluctuates infinitely often by a small percentage — as occurs in many splitting methods — the weights must shrink to zero.

Through martingale arguments, we establish that a splitting method provides asymptotically consistent estimates if and only if the sum of weights is fixed to one at all time steps. Moreover, under mild conditions, we show that a splitting method maintaining a constant sum of weights must be a WE method. Thus, we conclude that WE is the unique splitting method that provides asymptotically consistent estimates as $T \to \infty$.

1.2. *Variance bounds for weighted ensemble.* Our second contribution is to compare the accuracy of MCMC and WE estimates by considering the asymptotic variance as $T \to \infty$. For MCMC, there is a Central Limit Theorem that ensures the convergence in distribution

$$(1.3) \qquad \sqrt{T} \left( \frac{1}{T} \sum_{t=0}^{T-1} f(X_t) - \mu(f) \right) \xrightarrow{\mathcal{D}} \mathcal{N} \left( 0, \mu \left( v_f^2 \right) \right),$$

whenever the function $f$ is bounded and the MCMC sampler is geometrically ergodic [26]. In this Central Limit Theorem result, the asymptotic variance $\mu \left( v_f^2 \right)$ provides a quantitative measure of MCMC's accuracy and determines the simulation time that is needed to obtain accurate results. Here, the variance function $v_f^2$ is given explicitly by

$$(1.4) \qquad v_f = \sqrt{K h_f^2 - (K h_f)^2}, \quad h_f = \sum_{t=0}^{\infty} K^t (f - \mu(f))$$

where $K$ is the MCMC sampler's transition kernel [31, ch. 17].

Our work contributes new asymptotic variance bounds for WE that enable a comparison between the WE and MCMC. For a WE method with $N$ particles, we establish the lower bound

$$(1.5) \qquad \liminf_{T \to \infty} T \operatorname{Var} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f \left( \xi_t^i \right) \right] \geq \frac{\mu (v_f)^2}{N}.$$

Additionally, we prove the prefactor $\mu(v_f)^2$ is as sharp as possible. For any $\epsilon > 0$, we construct a WE scheme whose asymptotic variance satisfies

$$(1.6) \qquad \limsup_{T \to \infty} T \operatorname{Var} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f \left( \xi_t^i \right) \right] \leq (1 + \epsilon) \frac{\mu (v_f)^2}{N},$$

whenever the number of particles $N$ is sufficiently large. We prove these bounds for general unbounded functions $f$ under suitable integrability conditions.

Our asymptotic variance lower bound (1.5) is a fundamental result that restricts WE's behavior in all parameter regimes. This result is especially notable because many previous variance bounds for splitting schemes were derived in the mean field limit [15] in which the number of particles is large and the aggregate behavior becomes highly predictable. However, in our proof of the lower bound, we avoid reliance on the mean field limit; rather, we use direct variance manipulations and the assumption of geometric ergodicity to obtain a result that is valid for any number of particles $N$. In contrast, our proof that the optimal asymptotic variance can be approached from above does rely on the mean field limit as $N \to \infty$.

For large $T$ values, our results make it possible to quantify the maximal variance reduction of WE over MCMC in terms of an optimal improvement factor (OIF):

$$\text{(1.7)} \qquad \text{OIF} \equiv \frac{\mu(v_f^2)}{\mu(v_f)^2}.$$

When the OIF is large, as in many rare event probability estimation problems, our results guarantee the existence of a WE scheme that greatly increases efficiency compared to MCMC.

1.3. *Examples of weighted ensemble's efficiency.* In numerical examples, we demonstrates WE's usefulness for estimating rare event probabilities. These examples reveal that WE can provide dramatic benefits over MCMC, improving MCMC's variance by many orders of magnitude. Indeed, Figure 2 reveals a variance reduction of four orders of magnitude when calculating rare probabilities involving the Ising model (see Section 5.3 for details).
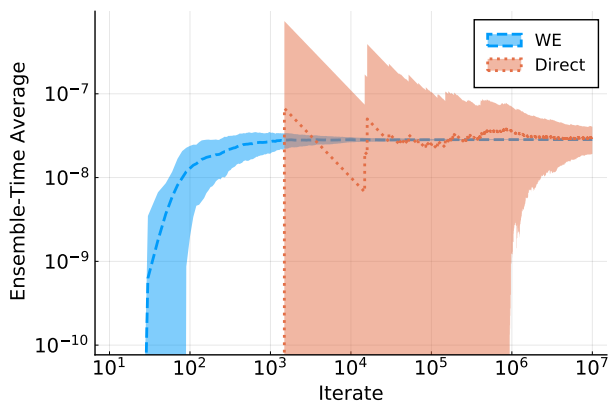


FIG 2. *Estimated probability of high magnetization in the high-temperature Ising model, as described in Section 5.3. The shaded region reflects one sample standard deviation, measured using* 100 *independent runs.*

These experiments add to the literature that demonstrates major efficiency gains by switching from MCMC to WE (e.g., [32]). These experiments also highlight a surprising consequence of our analysis. Typically, we would expect a Markov chain-based sampler to be less efficient than a sampler that directly draws independent samples from $\mu$, because the Markov chain-based samples are positively correlated [38]. However, in our experiments, WE produces estimates that are more efficient than could possibly be produced by an independence sampler. In conclusion, our work demonstrates how temporal correlations can be a strength, instead of a weakness, in rare event sampling.

1.4. *Outline for the paper.* The paper is organized as follows. Ergodicity theory is presented in Section 2, variance bounds are in Section 3, mathematical proofs are in Section 4, numerical experiments are in Section 5, and the conclusions follow in Section 6.

**2. Ergodicity theory for splitting schemes.** In this section, we define a splitting method and a weighted ensemble (WE) method. Then, we present our results proving that WE is the only splitting method that provides asymptotically consistent estimates as $T \to \infty$. Throughout the analysis, we use $\|f\| = \sup_x |f(x)|$ to denote the supremum norm on functions and $\|\mu\| = \sup_{\|f\| \leq 1} |\mu(f)|$ to denote the total variation norm on measures. We defer the technical proofs to Section 4.

2.1. *Definitions of splitting and weighted ensemble.* A splitting method [29, 36, 35] is a Monte Carlo method that alternates between a splitting step and an evolution step as follows.

ALGORITHM 2.1 (Splitting method).
First, independently sample particles $\xi_0^1, \ldots, \xi_0^{N_0}$ from a distribution $\mu_0$ and set $w_0^i = 1/N_0$ for $i = 1, \ldots, N_0$. Then, apply a splitting step and an evolution step at each time $t = 0, 1, \ldots$

1. Given particles and weights $\left(\xi_t^i, w_t^i\right)_{1 \leq i \leq N_t}$, apply the following splitting step.
   a. Select the mean number of children $C_t^i > 0$ for each particle $\xi_t^i$.
   b. Select the actual number of children $N_t^i \geq 0$ for each particle $\xi_t^i$, making sure that $N_t^i$ is a nonnegative integer with mean $C_t^i$.
   c. Split each particle $\xi_t^i$ into $N_t^i$ copies.
   d. Assign the children of $\xi_t^i$ uniform weights $w_t^i/C_t^i$.
2. Given particles and weights $\left(\hat{\xi}_t^i, \hat{w}_t^i\right)_{1 \leq i \leq N_{t+1}}$, apply the following evolution step.
   a. Evolve each particle $\hat{\xi}_t^i$ to a new state $\xi_{t+1}^i$ according to the transition kernel $K$.
   b. Assign each particle $\xi_{t+1}^i$ a weight $w_{t+1}^i = \hat{w}_t^i$.

A splitting method is highly general, since there are many possible strategies for choosing the numbers $C_t^i$ and $N_t^i$ during the splitting step. However, the main rule specified in Algorithm 2.1 is that children of $\xi_t^i$ receive uniform weights $w_t^i/C_t^i$. This rule ensures that the weights of the children of a $\xi_t^i$ sum up to the weight $w_t^i$ in expectation. To our knowledge, all the most popular splitting schemes are consistent with this rule, given appropriate definitions of $N_t^i$ and $C_t^i$. For example, in the original WE method of Huber and Kim [24], the $C_t^i$ are themselves random. Thus, a particle with weight $w_t^i = 1$ might randomly produce $C_t^i = N_t^i = 2$ copies with weights $1/2$ or $C_t^i = N_t^i = 3$ copies with weights $1/3$.

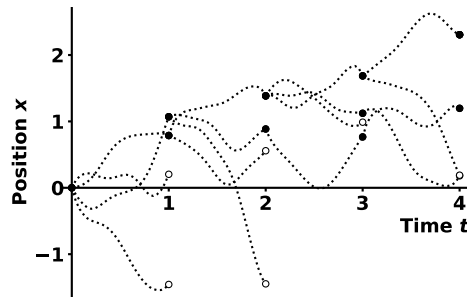For illustration, we show a typical splitting method in Figure 3 below.



FIG 3. *Splitting is used to sample rare, high values of the position $x$. White circles indicate that samples are killed. Black circles indicate that samples are preserved and possibly copied.*

A WE method [24, 13] is a particular type of splitting method that imposes more structure during the splitting step. In WE, we first divide particles into bins (more precisely, we divide particle indices $1, 2, \ldots, N_t$ into bins). Then, we use splitting and killing to adjust the populations in the bins while exactly preserving the bins' statistical weights.

ALGORITHM 2.2 (Weighted ensemble).
Apply a splitting method, and at each time $t \geq 0$ perform the following splitting step:

a. Partition the indices $1 \leq i \leq N_t$ into bins $u_1, u_2, \ldots,$ and set $w_t(u) = \sum_{i \in u} w_t^i$.
b. Select the desired number of children $N_t(u) \geq 1$ for each bin $u$.

c. Select the number of children $N_t^i \geq 0$ for each particle $\xi_t^i$ in each bin $u$, making sure $N_t^i$ is a nonnegative integer with expectation $C_t^i = N_t(u) w_t^i / w_t(u)$ and $\sum_{i \in u} N_t^i = N_t(u)$.
d. Split each particle $\xi_t^i$ into $N_t^i$ copies.
e. Assign the children with parents in bin $u$ uniform weights $w_t(u)/N_t(u)$.

Our definition of weighted ensemble allows for an arbitrary choice of bins that may change at each time step $t \geq 0$. However, in our analysis, we focus on the simple case where bins correspond to fixed regions that divide up the state space, as was shown in Figure 1. Past work exploring optimal bin design strategies includes [4, 10, 11, 39].

To complete our introduction to splitting and WE methods, we describe common approaches for selecting a population of children particles $\hat{\xi}_t^1, \ldots, \hat{\xi}_t^{N_{t+1}}$ from a population of parent particles $\xi_t^1, \ldots, \xi_t^{N_t}$. We assume that the WE user has already determined that each particle $\xi_t^i$ should produce $C_t^i > 0$ children particles on average, and we describe several *resampling schemes* [40] for determining the precise number of children $N_t^i$ for each particle $\xi_t^i$.

The simplest resampling scheme is *multinomial resampling*, which we describe below.

DEFINITION 2.1. In multinomial resampling [15], we independently sample children particles $\hat{\xi}_t^1, \ldots, \hat{\xi}_t^{N_{t+1}}$ from locations $(\xi_t^i)_{1 \leq i \leq N_t}$ with probabilities proportional to $(C_t^i)_{1 \leq i \leq N_t}$. Thus, the numbers $(N_t^i)_{1 \leq i \leq N_t}$ are jointly distributed according to

$$(2.1) \qquad \left( N_t^1, \ldots, N_t^{N_t} \right) \sim \text{Multi} \left( N_t, \frac{C_t^1}{N_t}, \ldots, \frac{C_t^{N_t}}{N_t} \right).$$

Multinomial resampling is used in many splitting schemes [16], but it cannot be used in weighted ensemble since it would violate the requirement of placing exactly $N_t(u)$ children particles in each bin $u$. However, a related approach called *binned multinomial resampling* can be used with WE instead.

DEFINITION 2.2. In binned multinomial resampling, we iterate over the bins and apply multinomial resampling within each bin. Thus, if the particles in bin $u$ are $\xi_t^{i_1}, \ldots, \xi_t^{i_m}$, the numbers $N_t^{i_1}, \ldots, N_t^{i_m}$ are jointly distributed according to

$$(2.2) \qquad \left( N_t^{i_1}, \ldots, N_t^{i_m} \right) \sim \text{Multi} \left( N_t(u), \frac{w_t^{i_1}}{w_t(u)}, \ldots, \frac{w_t^{i_m}}{w_t(u)} \right).$$

In multinomial resampling and binned multinomial resampling, we observe that children particles $\hat{\xi}_t^1, \ldots, \hat{\xi}_t^{N_{t+1}}$ are independently sampled given auxiliary information including the locations of the parents and the mean number of children for each parent. In general, we can consider other resampling possible schemes that maintain this conditional independence property. We define the class of *conditionally independent resampling* schemes as follows.

DEFINITION 2.3. In a conditionally independent resampling scheme [40], given parent particles $\xi_t^1, \ldots, \xi_t^{N_t}$ with weights $w_t^1, \ldots, w_t^{N_t}$, we define a matrix $\boldsymbol{P} \in \mathbb{R}^{M \times N_t}$ where $M$ represents the maximum possible number of allowable children. Then, we iterate over $i = 1, 2, \ldots, M$. With probability $\boldsymbol{P}_{ij}$, we assign

$$(2.3) \qquad \hat{\xi}_t^i = \xi_t^j \quad \text{and} \quad \hat{w}_t^i = \frac{w_t^j}{C_t^j}.$$

With the remaining probability $1 - \sum_{j=1}^{N_t} \boldsymbol{P}_{ij}$, we do not assign $\hat{\xi}_t^i$ to any location at all and we set $\hat{w}_t^i = 0$. We remove particles with zero weights at the end of the resampling scheme.

Conditionally independent resampling schemes are a broad category that encompasses most procedures that are used in practice. In addition to multinomial resampling and binned multinomial resampling, this category includes Bernoulli resampling, multinomial residual resampling, stratified resampling, and stratified residual resampling [18, 40].

While the choice of resampling scheme can affect a splitting method's variance, here we do not provide a detailed comparison between different procedures. Rather, we emphasize broad results that hold for many different resampling schemes. In our analysis, we only make specific assumptions about the resampling scheme twice. First, when we establish that WE is the unique splitting method providing asymptotically convergent estimates, we assume a conditionally independent resampling scheme (see Proposition 2.2). Second, in our our demonstration that WE can approach the optimal asymptotic variance from above, our construction is based on binned multinomial resampling (see Lemma 3.1).

2.2. *Ergodicity of splitting schemes.* When we combine splitting with MCMC, we must carefully consider the long-time behavior of the splitting method's estimates. It is well-known that MCMC estimates must converge

$$(2.4) \qquad \frac{1}{T} \sum_{t=0}^{T-1} f\left(X_t\right) \overset{T\to\infty}{\to} \mu\left(f\right),$$

assuming $f$ is bounded and the dynamics are Harris ergodic [31, ch. 17]. Therefore, we ask whether splitting estimates also converge similarly to MCMC estimates. Specifically, we ask: do the estimates from a splitting method always satisfy

$$(2.5) \qquad \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right) \overset{T\to\infty}{\to} \mu\left(f\right),$$

and, if not, what assumptions guarantee the convergence in (2.5)?

To address these questions in a rigorous way, we first define the ergodicity conditions that will be considered in the analysis.

DEFINITION 2.4 (Ergodicity conditions).
Consider a $\psi$-irreducible, aperiodic transition kernel $K$ on a general state space $X$, and assume $K$ is invariant with respect to a distribution $\mu = \mu K$.

(i) The kernel $K$ is *Harris ergodic* if

$$(2.6) \qquad \left\| K^t\left(x, \cdot\right) - \mu \right\| \overset{t\to\infty}{\to} 0$$

for all $x \in X$.

(ii) The kernel $K$ is *geometrically ergodic* if

$$(2.7) \qquad \sum_{t=0}^{\infty} r^t \left\| K^t\left(x, \cdot\right) - \mu \right\| < \infty$$

for fixed $r > 1$ and all $x \in X$.

(iii) The kernel $K$ is *V-uniformly ergodic* for a function $1 \le V \le \infty$ if $\mu\left(V\right) < \infty$ and

$$(2.8) \qquad \sup_{|g| \le V} \left| K^t g\left(x\right) - \mu\left(g\right) \right| \le R \rho^t V(x)$$

for fixed $R > 0$, fixed $\rho < 1$, and all $x \in X$.

Harris ergodicity is a comparatively weak condition that gives no control over the convergence rate in $\left\| K^t\left(x,\cdot\right) - \mu \right\| \overset{t\to\infty}{\to} 0$, whereas geometric ergodicity and $V$-uniform ergodicity are stronger conditions that specify an exponential convergence rate. While geometric ergodicity and $V$-uniform ergodicity are nearly equivalent, we exploit the slight difference between these conditions in our analysis. As explained in [31, ch.15], geometric ergodicity implies $V$-uniform ergodicity for a particular function $V$. Conversely, $V$-uniform ergodicity implies geometric ergodicity if we restrict the process to the absorbing set $\{V < \infty\}$.

In contrast to an MCMC method, our experiments reveal that a splitting method does not necessarily provide consistent estimates as $T \to \infty$, even when $K$ is geometrically ergodic and $f$ is bounded. Figure 4 shows an example of a splitting scheme that fails to provide consistent estimates. The sum of the weights approaches zero over long timescales, which causes estimates to converge to zero also.
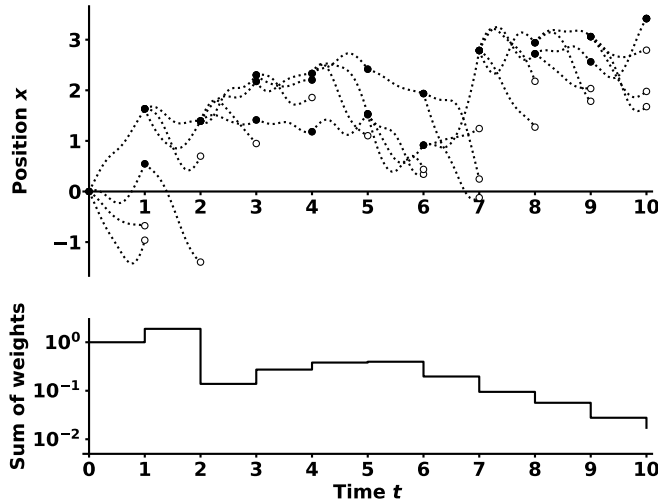


FIG 4. *Weights converge to zero as splitting and killing are repeatedly applied.*

This problem of shrinking weights was first pointed out by Aristoff in [2], and here we provide a full mathematical explanation. Since the sum of the weights is a martingale, small relative fluctuations in the sum of the weights build up over time. Moreover, the small fluctuations lead to major consequences, as we demonstrate in the proposition below.

PROPOSITION 2.1.    *In a splitting method, suppose there exists $\epsilon > 0$ such that the event*

$$(2.9) \qquad \left| \frac{\sum_{i=1}^{N_{t+1}} w_{t+1}^i}{\sum_{i=1}^{N_t} w_t^i} - 1 \right| > \epsilon$$

*occurs infinitely often with probability one. Then, almost surely, $\sum_{i=1}^{N_t} w_t^i \to 0$ as $t \to \infty$.*

As a consequence of Proposition 2.1, the only way to avoid shrinking weights is to asymptotically eliminate all small fluctuations as $T \to \infty$. This pressing need to control the sum of the weights leads us to consider the possibility of simply fixing the sum of weights to be one, as naturally occurs in WE. In the next theorem, we verify that a splitting scheme provides asymptotically consistent estimates if and only if $\sum_{i=1}^{N_t} w_t^i = 1$ at all times $t \geq 0$.

THEOREM 2.1. *Consider a splitting method with a $V$-uniformly ergodic kernel $K$ and assume $\mu_0 \{V < \infty\} = 1$. Then, the following three conditions are equivalent:*

(i) *The time average of the sum of the weights converges in probability to one:*

$$(2.10) \qquad P\left\{\left|\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N_t} w_t^i - 1\right| > \epsilon\right\} \stackrel{T\to\infty}{\to} 0, \quad \epsilon > 0.$$

(ii) *With probability one, the weights satisfy $\sum_{i=1}^{N} w_t^i = 1$ at all times $t \geq 0$.*
(iii) *Whenever $\left\|f^2/V\right\| < \infty$, the estimates of $\mu(f)$ converge with probability one:*

$$(2.11) \qquad P\left\{\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) \stackrel{T\to\infty}{\to} \mu(f)\right\} = 1.$$

Next, we prove that the *only* splitting method capable of maintaining a constant sum of weights is WE, assuming a conditionally independent resampling scheme is used.

PROPOSITION 2.2. *If a splitting method with a conditionally independent resampling scheme satisfies $\sum_{i=1}^{N_t} w_t^i = 1$ at all times $t \geq 0$, the splitting method is a weighted ensemble method with a particular choice of bins.*

In summary, we have proved under mild conditions that WE is the only splitting method that provides asymptotically consistent estimates as $T \to \infty$.

REMARK 2.1. In this section, we have analyzed a splitting method's estimates

$$(2.12) \qquad \mu(f) \approx \frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right).$$

Yet, we might also consider a splitting method's normalized estimates

$$(2.13) \qquad \mu(f) \approx \frac{1}{T}\sum_{t=0}^{T-1}\frac{\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right)}{\sum_{i=1}^{N_t} w_t^i}.$$

In WE, the normalized and unnormalized estimates are the same. However, in splitting methods other than WE, the normalized and unnormalized estimates differ. Past analyses of splitting methods [15, 23] showed that the normalized estimates typically converge as $T \to \infty$ with an asymptotic bias of size $\mathcal{O}\left(\frac{1}{N}\right)$, preventing the possibility of consistent estimation. WE presents the only known exception to this trend, since the asymptotic bias is zero.

**3. The bias and variance of weighted ensemble estimates.** In this section, our broad goal is to determine whether WE can produce more accurate estimates than MCMC. We first analyze the bias and then analyze the variance of WE estimates. Throughout the section, we fix the number of particles to be exactly $N$ at all time steps, and we analyze $N$ as a key control parameter influencing WE's efficiency.

3.1. *Bias of weighted ensemble estimates.* As our first result, we find that WE adds no additional bias compared to MCMC. Rather, as was originally shown in [41], WE estimates have the same expectation as estimates from a standard MCMC sampler.

PROPOSITION 3.1. *Consider a WE scheme with a Harris ergodic kernel $K$, and assume $f$ is bounded. WE estimates for $\mu(f)$ have the following bias properties:*

1. *For any $T \geq 0$, the WE estimate*

$$(3.1) \qquad \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)$$

*has the same expectation as a trajectory average*

$$(3.2) \qquad \frac{1}{T} \sum_{t=0}^{T-1} f\left(X_t\right),$$

*where $X_t$ is a Markov chain with transition kernel $K$ and initial distribution $\mu_0$.*

2. *The WE estimates for $\mu(f)$ are asymptotically unbiased in the limit as $T \to \infty$:*

$$(3.3) \qquad \mathrm{E}\left[\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)\right] \overset{T \to \infty}{\Rightarrow} \mu(f).$$

3. *With a burn-in period of length $\tau$, the WE estimates for $\mu(f)$ satisfy*

$$(3.4) \qquad \mathrm{E}\left[\frac{1}{T} \sum_{t=\tau}^{\tau+T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)\right] \overset{\tau \to \infty}{\Rightarrow} \mu(f).$$

In the limit as $T \to \infty$, Proposition 3.1 shows that WE is asymptotically unbiased. However, we may worry about bias in the pre-asymptotic regime. To reduce bias, therefore, we can run the WE algorithm for an extra $\tau$ time steps and use the estimate

$$(3.5) \qquad \mu(f) \approx \frac{1}{T} \sum_{t=\tau}^{\tau+T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right).$$

Proposition 3.1 verifies that incorporating a "burn-in period" of length $\tau$ has a beneficial impact. As $\tau \to \infty$, the bias in WE estimates vanishes completely.

3.2. *Variance of weighted ensemble estimates.* Now that we have considered bias, our next step is evaluating the variance of WE estimates. To provide simple formulas for the WE variance, we fix a function $f$ and define the associated conditional expectation function

$$(3.6) \qquad h_f(x) = \sum_{t=0}^{\infty} \left(K^t f(x) - \mu(f)\right)$$

and the function

$$(3.7) \qquad v_f(x) = \sqrt{K h_f^2(x) - \left(K h_f(x)\right)^2}.$$

The function $v_f^2$ is commonly used in the Markov chain literature to express the asymptotic variance of trajectory averages involving $f$. Here, we build on this literature by using $v_f$ to also compare MCMC and WE variances. Our main result is the following theorem, which establishes the best possible asymptotic variance for WE estimates.

THEOREM 3.1. *Consider a WE scheme with a kernel $K$ that is geometrically ergodic and $V$-uniformly ergodic, and assume $\left\|f^2/V\right\| < \infty$. WE estimates for $\mu(f)$ have the following variance properties:*

1. *The variance of WE estimates is bounded from below by*

$$(3.8) \qquad \liminf_{T \to \infty} T \operatorname{Var}\left[\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)\right] \geq \frac{\mu\left(v_f\right)^2}{N}.$$

2. *For any $\epsilon > 0$, there is a particular WE scheme requiring a sufficiently large number of particles $N$ that satisfies*

$$(3.9) \qquad \limsup_{T \to \infty} T \operatorname{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)\right] \le (1+\epsilon)\frac{\mu\left(v_f\right)^2}{N}.$$

Theorem 3.1 opens up the possibility for a quantitative comparison between MCMC's and WE's efficiencies. For MCMC, we can either run a single Markov chain for $NT$ time steps, or we can run $N$ independent Markov chains for $T$ time. Both approaches share a similar computational cost, and both approaches yield an estimate of $\mu(f)$ whose variance is nearly $\mu\left(v_f^2\right)/NT$. In contrast, by running WE for $T$ time steps with $N$ particles, it may be possible to achieve a much lower variance. With the optimal design parameters and with a sufficiently large number of particles $N$, WE produces an estimate of $\mu(f)$ whose variance is nearly $\mu(v_f)^2/NT$.

We can quantify the efficiency benefits of WE over MCMC by means of the optimal improvement factor (OIF) that was previously discussed in the introduction section:

$$(3.10) \qquad \mathrm{OIF} \equiv \frac{\mu(v_f^2)}{\mu(v_f)^2}.$$

If this factor is one, then WE cannot improve MCMC's variance at all — a situation that occurs, for example, if the kernel $K$ is an independence sampler. However, in rare event probability estimation, the OIF is typically multiple multiple orders of magnitude, demonstrating major potential for WE to reduce MCMC's variance. In Section 5, we explicitly calculate the OIF for several examples.

We close this section by discussing the key lemma that allows us to optimize WE's variance and prove the sharpness result (3.9).

LEMMA 3.1. *Consider a WE scheme with a $V$-uniformly ergodic kernel $K$. Assume binned multinomial resampling is used, $\mu_0(V) < \infty$, and $\left\|f^2/V\right\| < \infty$. Then, as $T \to \infty$, WE estimates for $\mu(f)$ satisfy*

$$(3.11) \quad \operatorname{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)\right]$$

$$(3.12) \quad = \frac{1}{T^2}\sum_{t=0}^{T-2} \mathrm{E}\left[\sum_u \frac{w_t\left(u\right)^2}{N_t\left(u\right)}\left(\operatorname{Var}_{\eta_t^u}\left[Kh_f\right] + \operatorname{Var}_{\eta_t^u}\left[v_f\right] + \eta_t^u\left(v_f\right)^2\right)\right] + \mathcal{O}\left(\frac{1}{T^2}\right),$$

*where $\eta_t^u = \frac{1}{w_t(u)}\sum_{i\in u} w_t^i \delta_{\xi_t^i}$ denotes the empirical distribution of particles in bin $u$.*

This lemma decomposes WE's asymptotic variance into separate contributions from the different bins. The lemma reveals how these bins can be optimized to minimize WE's variance. In our proof of (3.9), we choose the bins and bin allocations in the following way:

1. We first define a large number of spatial bins in the $Kh_f$ and $v_f$ coordinates, ensuring that most of the terms $\operatorname{Var}_{\eta_t^u}\left[Kh_f\right] + \operatorname{Var}_{\eta_t^u}\left[v_f\right]$ in the variance decomposition are small.
2. We next minimize the $w_t\left(u\right)^2 \eta_t^u\left(v_f\right)^2/N_t\left(u\right)$ terms in the variance decomposition by allocating particles to bins according to the rule

$$(3.13) \qquad \frac{N_t\left(u\right)}{N} \approx \frac{w_t\left(u\right)\eta_t^u\left(v_f\right)}{\sum_{u'} w_t\left(u'\right)\eta_t^{u'}\left(v_f\right)}.$$

3. As $\epsilon \to 0$, we increase the number of particles and bins to ensure that WE's variance lies within a factor of $1 + \epsilon$ of the optimal variance $\mu (v_f)^2 / NT$.

While this optimization strategy is convenient for proving the variance bound (3.9), it would be difficult to carry out this strategy in WE applications. The main problem is that functions $Kh_f$ and $v_f$ are typically unknown. As a more practical alternative, therefore, Aristoff and Zuckerman [4] have developed an optimization approach for WE that uses coarse-grained approximations of the functions $Kh_f$ and $v_f$. We apply this optimization approach in all the numerical examples in Section 5.

**4. Mathematical proofs.** Here, we prove our theoretical results concerning the bias, convergence, and variance of WE estimates.

4.1. *Bias.* We first examine the bias of a splitting method's estimates. As a central analysis tool, we consider a filtration of $\sigma$-algebras $\mathcal{F}_0 \subseteq \hat{\mathcal{F}}_0 \subseteq \mathcal{F}_1 \subseteq \hat{\mathcal{F}}_1 \subseteq \cdots$ that satisfy the following assumptions:

ASSUMPTIONS 4.1.

(i) The variables $\left( \xi_t^i, w_t^i, C_t^i \right)_{1 \le i \le N_t}$ are measurable with respect to $\mathcal{F}_t$.

(ii) Conditional on $\mathcal{F}_t$, the copy numbers $N_t^1, \ldots, N_t^{N_t}$ each have mean $\mathrm{E}\left[ N_t^i \middle| \mathcal{F}_t \right] = C_t^i$.

(iii) The variables $\left( \hat{\xi}_t^i, \hat{w}_t^i \right)_{1 \le i \le N_{t+1}}$ are measurable with respect to $\hat{\mathcal{F}}_t$.

(iv) Conditional on $\hat{\mathcal{F}}_t$, the particles $\xi_{t+1}^1, \ldots, \xi_{t+1}^{N_{t+1}}$ are independent with $\mathrm{Law}\left( \xi_{t+1}^i \middle| \hat{\mathcal{F}}_t \right) = K\left( \hat{\xi}_t^i, \cdot \right)$.

The filtration $\mathcal{F}_0 \subseteq \hat{\mathcal{F}}_0 \subseteq \mathcal{F}_1 \subseteq \hat{\mathcal{F}}_1 \subseteq \cdots$ has a natural interpretation in terms of the information that is available at each step of the splitting method. $\mathcal{F}_t$ contains all the information available after the identity of the particles $\xi_t^1, \ldots, \xi_t^{N_t}$ is revealed and before the identities of the children particles $\hat{\xi}_t^1, \ldots, \hat{\xi}_t^{N_{t+1}}$ are revealed. $\hat{\mathcal{F}}_t$ contains all the information in $\mathcal{F}_t$ and also the identities of the children particles $\hat{\xi}_t^1, \ldots, \hat{\xi}_t^{N_{t+1}}$.

The $\sigma$-algebras are useful because they reveal a rich martingale structure that underlies splitting schemes, which was originally exploited by Del Moral in [15]. We introduce this martingale structure in the following lemma:

LEMMA 4.1. *Fix a time $T \ge 0$ and a function $f$ with $\mu_0 K^T |f| < \infty$. Define*

$$(4.1) \quad M_t = \mathrm{E}\left[ \sum_{i=1}^{N_T} w_T^i f\left( \xi_T^i \right) \middle| \mathcal{F}_t \right], \quad \hat{M}_t = \mathrm{E}\left[ \sum_{i=1}^{N_T} w_T^i f\left( \xi_T^i \right) \middle| \hat{\mathcal{F}}_t \right], \quad 0 \le t \le T-1.$$

*Then, $M_0, \hat{M}_0, \ldots, M_{T-1}, \hat{M}_{T-1}$ is a martingale that satisfies*

$$(4.2) \quad M_t = \sum_{i=1}^{N_t} w_t^i \left( K^{T-t} f \right)\left( \xi_t^i \right), \quad \hat{M}_t = \sum_{i=1}^{N_{t+1}} \hat{w}_t^i \left( K^{T-t} f \right)\left( \hat{\xi}_t^i \right), \quad 0 \le t \le T-1.$$

PROOF. Set $M_T = \sum_{i=1}^{N_T} w_T^i f\left( \xi_T^i \right)$ and assume for some $0 \le t \le T-1$ the representation $M_{t+1} = \sum_{i=1}^{N_{t+1}} w_{t+1}^i \left( K^{T-t-1} f \right)\left( \xi_{t+1}^i \right)$ is valid. Then, using Assumption 4.1 (iv),

$$(4.3) \quad \hat{M}_t = \mathrm{E}\left[ M_{t+1} \middle| \hat{\mathcal{F}}_t \right]$$

$$(4.4) \qquad = \sum_{i=1}^{N_{t+1}} \mathrm{E}\left[ w_{t+1}^i \left( K^{T-t-1} f \right) \left( \xi_{t+1}^i \right) \middle| \mathcal{F}_t \right]$$

$$(4.5) \qquad = \sum_{i=1}^{N_{t+1}} \hat{w}_t^i \left( K^{T-t} f \right) \left( \hat{\xi}_t^i \right).$$

Using Assumption 4.1 (ii) and the fact that children of $\xi_t^i$ receive weights $w_t^i/C_t^i$,

$$(4.6) \qquad M_t = \mathrm{E}\left[ \hat{M}_t \middle| \mathcal{F}_t \right]$$

$$(4.7) \qquad = \mathrm{E}\left[ \sum_{i=1}^{N_{t+1}} \hat{w}_t^i \left( K^{T-t} f \right) \left( \hat{\xi}_t^i \right) \middle| \mathcal{F}_t \right]$$

$$(4.8) \qquad = \mathrm{E}\left[ \sum_{i=1}^{N_t} N_t^i \frac{w_t^i}{C_t^i} \left( K^{T-t} f \right) \left( \xi_t^i \right) \middle| \mathcal{F}_t \right]$$

$$(4.9) \qquad = \sum_{i=1}^{N_t} w_t^i \left( K^{T-t} f \right) \left( \xi_t^i \right).$$

$\square$

Lemma 4.1 allows us to prove the following generalization of Proposition 3.1, which simultaneously establishes bias properties for all splitting methods and WE methods.

PROPOSITION 4.1. *Consider a splitting method with a Harris ergodic kernel $K$, and assume $f$ is bounded. The estimates for $\mu(f)$ have the following bias properties:*

1. *For any $\tau \geq 0$ and any $T \geq 0$, the estimate*

$$(4.10) \qquad \frac{1}{T} \sum_{t=\tau}^{\tau+T-1} \sum_{i=1}^{N_t} w_t^i f \left( \xi_t^i \right)$$

   *has the same expectation as the trajectory average*

$$(4.11) \qquad \frac{1}{T} \sum_{t=\tau}^{\tau+T-1} f \left( X_t \right),$$

   *where $X_t$ is a Markov chain with transition kernel $K$ and initial distribution $\mu_0$.*
2. *The estimates for $\mu(f)$ are asymptotically unbiased in the limit as $\tau + T \to \infty$:*

$$(4.12) \qquad \mathrm{E}\left[ \frac{1}{T} \sum_{t=\tau}^{\tau+T-1} \sum_{i=1}^{N_t} w_t^i f \left( \xi_t^i \right) \right] \stackrel{\tau+T\to\infty}{\longrightarrow} \mu(f).$$

PROOF. Using Lemma 4.1, we calculate

$$(4.13) \qquad \mathrm{E}\left[ \frac{1}{T} \sum_{t=\tau}^{\tau+T-1} \sum_{i=1}^{N_t} w_t^i f \left( \xi_t^i \right) \right] = \frac{1}{T} \sum_{t=\tau}^{T+\tau-1} \mu_0 K^t f.$$

As a consequence of Harris ergodicity, we have the convergence $\left\| \mu_0 K^t - \mu \right\| \stackrel{t\to\infty}{\longrightarrow} 0$ [31, ch. 13]. Sending $\tau + T \to \infty$, we verify

$$(4.14) \qquad \left| \mathrm{E}\left[ \frac{1}{T} \sum_{t=\tau}^{\tau+T-1} \sum_{i=1}^{N_t} w_t^i f \left( \xi_t^i \right) \right] - \mu(f) \right| \leq \frac{\|f\|}{T} \sum_{t=\tau}^{\tau+T-1} \left\| \mu_0 K^t - \mu \right\| \to 0.$$

□

4.2. *Convergence.*  In this section, we prove that a splitting method provides asymptotically consistent estimates if and only if the sum of the weights is almost surely one. To prove this result, we observe that the splitting method defined in Algorithm 2.1 ensures that the sum of the weights has expected value one at all times $t \geq 0$. Moreover, the sum of the weights $\sum_{i=1}^{N_t} w_t^i$ is a nonnegative martingale, and a nonnegative martingale must converge with probability one as $t \to \infty$ [28]. This observation immediately verifies the result in Proposition 2.1. To prove Theorem 2.1, we also need the following lemma:

LEMMA 4.2.   *If $K$ is $V$-uniformly ergodic, then $K$ is also $\sqrt{V}$-uniformly ergodic.*

PROOF.  By Jensen's inequality, for any positive measure $\eta$,

$$\tag{4.15} \sup_{|g| \leq \sqrt{V}} \eta\left(|g|\right) \leq \|\eta\| \sup_{g^2 \leq V} \frac{\eta}{\|\eta\|}\left(|g|\right)$$

$$\tag{4.16} \leq \|\eta\| \sup_{g^2 \leq V} \sqrt{\frac{\eta}{\|\eta\|}\left(g^2\right)}$$

$$\tag{4.17} = \sqrt{\|\eta\|} \sqrt{\sup_{|g| \leq V} \eta\left(|g|\right)}.$$

Taking $\eta = \left|K^t\left(x, \cdot\right) - \mu\right|$ and applying $V$-uniform ergodicity gives the desired result.   □

PROOF OF THEOREM 2.1.  First, we observe that (iii) implies (i).

Next, we show that (i) implies (ii). Part (i) indicates the convergence in probability $\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N_t} w_t^i \overset{T \to \infty}{\to} 1$. Since $\sum_{i=1}^{N_t} w_t^i$ is a nonnegative martingale, $\sum_{i=1}^{N_t} w_t^i$ converges almost surely to a random variable $W_\infty$ as $t \to \infty$. Hence, we must have $W_\infty = 1$. Next, for fixed $t \geq 0$, Fatou's lemma implies

$$\tag{4.18} \sum_{i=1}^{N_t} w_t^i = \liminf_{T \to \infty} \mathrm{E}\left[\left.\sum_{i=1}^{N_T} w_T^i \right| \mathcal{F}_t\right] \geq \mathrm{E}\left[\left.\liminf_{T \to \infty} \sum_{i=1}^{N_T} w_T^i \right| \mathcal{F}_t\right] = \mathrm{E}\left[W_\infty | \mathcal{F}_t\right] = 1.$$

Since $\mathrm{E}\left[\sum_{i=1}^{N_t} w_t^i\right] = 1$ and $\sum_{i=1}^{N_t} w_t^i \geq 1$, we conclude that $\sum_{i=1}^{N_t} w_t^i = 1$ with probability one. Since $t \geq 0$ is arbitrary, we have verified (ii).

Last of all, we prove that (ii) implies (iii). We assume without loss of generality $f \geq 0$, and we show that almost surely

$$\tag{4.19} \mathrm{P}\left\{\left.\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) \overset{T \to \infty}{\to} \int \mu\left(dx\right) f\left(x\right) \right| \mathcal{F}_0\right\} = 1.$$

We fix $T \geq 0$ and compute the conditional variance

$$\tag{4.20} \mathrm{Var}\left[\left.\sum_{t=0}^{T-1} \sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) \right| \mathcal{F}_0\right] = \sum_{s,t=0}^{T-1} \mathrm{Cov}\left[\left.\sum_{i=1}^{N_s} w_s^i f\left(\xi_s^i\right), \sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) \right| \mathcal{F}_0\right].$$

For $s \leq t$, the conditional covariance terms satisfy

$$\tag{4.21} \mathrm{Cov}\left[\left.\sum_{i=1}^{N_s} w_s^i f\left(\xi_s^i\right), \sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) \right| \mathcal{F}_0\right]$$

$$(4.22) \quad = \mathrm{Cov}\left[\sum_{i=1}^{N_s} w_s^i f\left(\xi_s^i\right), \sum_{i=1}^{N_s} w_s^i K^{t-s} f\left(\xi_s^i\right) \,\middle|\, \mathcal{F}_0\right]$$

$$(4.23) \quad \leq \mathrm{Var}\left[\sum_{i=1}^{N_s} w_s^i f\left(\xi_s^i\right) \,\middle|\, \mathcal{F}_0\right]^{1/2} \mathrm{Var}\left[\sum_{i=1}^{N_s} w_s^i K^{t-s} f\left(\xi_s^i\right) \,\middle|\, \mathcal{F}_0\right]^{1/2}.$$

Using the fact that $\sum_{i=1}^{N_s} w_s^i = 1$, we calculate

$$(4.24) \quad \mathrm{Var}\left[\sum_{i=1}^{N_s} w_s^i K^{t-s} f\left(\xi_s^i\right) \,\middle|\, \mathcal{F}_0\right]$$

$$(4.25) \quad \leq \mathrm{E}\left[\left|\sum_{i=1}^{N_s} w_s^i \left(K^{t-s} f - \mu\left(f\right)\right)\left(\xi_s^i\right)\right|^2 \,\middle|\, \mathcal{F}_0\right]$$

$$(4.26) \quad \leq \mathrm{E}\left[\sum_{i=1}^{N_s} w_s^i \left(K^{t-s} f - \mu\left(f\right)\right)^2 \left(\xi_s^i\right) \,\middle|\, \mathcal{F}_0\right]$$

$$(4.27) \quad = \frac{1}{N_0}\sum_{i=1}^{N_0} K^s \left(\left(K^{t-s} f - \mu\left(f\right)\right)^2\right)\left(\xi_0^i\right)$$

Using the $\sqrt{V}$-uniform ergodicity and $V$-uniform ergodicity of $K$, the last term is size $\mathcal{O}\left(r^{-(t-s)}\right)$ for a fixed constant $r > 1$, and we obtain a bound of the form

$$(4.28) \quad \mathrm{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) \,\middle|\, \mathcal{F}_0\right] \leq \frac{C\left\|f^2/V\right\|}{T},$$

where $C$ is independent of $T$ and $f$. Since the conditional variance terms are summable for $T = 1, 4, 9, \ldots$, the WE estimates converge by a Borel-Cantelli argument, and we find

$$(4.29) \quad \lim_{T\to\infty}\frac{1}{T^2}\sum_{t=0}^{T^2-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) = \lim_{T\to\infty}\mathrm{E}\left[\frac{1}{T^2}\sum_{t=0}^{T^2-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right) \,\middle|\, \mathcal{F}_0\right] = \mu\left(f\right),$$

with conditional probability one. We can strengthen the almost sure convergence for $T = 1, 4, 9, \ldots$ to almost sure convergence for $T = 1, 2, 3, \ldots$ by noting that

$$(4.30) \quad \frac{T^2}{T^2+s}\left(\frac{1}{T^2}\sum_{t=0}^{T^2-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right)\right) \leq \frac{1}{T^2+s}\sum_{t=0}^{T^2+s-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right)$$

$$(4.31) \quad \leq \frac{(T+1)^2}{T^2+s}\left(\frac{1}{(T+1)^2}\sum_{t=0}^{(T+1)^2-1}\sum_{i=1}^{N_t} w_t^i f\left(\xi_t^i\right)\right)$$

whenever $T^2 \leq T^2 + s \leq (T+1)^2$. Hence, we verify equation (4.19), completing the proof. $\square$

Lastly, we verify that any conditionally independent resampling scheme that maintains a sum of weights equal to one is a WE scheme:

PROOF OF PROPOSITION 2.2. We condition on the matrix $\boldsymbol{P} \in \mathbb{R}^{M \times N_t}$ and on the locations and weights of the parents. Before removing the particles with zero weights, the weights $\hat{w}_t^1, \ldots, \hat{w}_t^M$ are independent. Since $\sum_{i=1}^M \hat{w}_t^i = 1$, we find

$$(4.32) \qquad 0 = \operatorname{Var}\left[\sum_{i=1}^M \hat{w}_t^i\right] = \sum_{i=1}^M \operatorname{Var}\left[\hat{w}_t^i\right].$$

and each weight $\hat{w}_t^i$ is constant with probability one. Hence, we can define bins $u_c$ consisting of all the parents whose children receive weights $\hat{w}_t^i = c$. There is a fixed number of children per bin and all the children receive the same weight, so the splitting method is a WE method. $\square$

4.3. *Variance.* In this final subsection of technical results, we bound the variance of WE estimates. Our main approach, following the analysis developed by Del Moral [15], is to decompose the variance of WE estimates as a sum of squared martingale differences and then manipulate the martingale difference terms to obtain sharp error bounds.

The martingale we use is slightly different from the one described in Lemma 4.1, since we need to account for the time-averaging that produces WE estimates. The following lemma introduces this martingale and gives an explicit formula for the martingale differences:

LEMMA 4.3. *Fix $T \geq 0$ and a function $f$ with $\mu_0 K^t |f| < \infty$ for $0 \leq t \leq T - 2$. Define*

$$(4.33) \qquad Y_t = \mathrm{E}\left[\frac{1}{T}\sum_{t=0}^{t-1}\sum_{i=1}^N w_T^i f\left(\xi_T^i\right)\middle|\mathcal{F}_t\right], \qquad \hat{Y}_t = \mathrm{E}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^N w_T^i f\left(\xi_T^i\right)\middle|\hat{\mathcal{F}}_t\right].$$

*Then, $Y_0, \hat{Y}_0, \ldots, Y_{T-1}, \hat{Y}_{T-2}$ is a martingale with martingale differences given by*

$$(4.34) \qquad \hat{Y}_t - Y_t = \frac{1}{T}\left[\sum_{i=1}^N \hat{w}_t^i K h_{t+1}^T\left(\hat{\xi}_t^i\right) - \sum_{i=1}^N w_t^i K h_{t+1}^T\left(\xi_t^i\right)\right],$$

$$(4.35) \qquad Y_{t+1} - \hat{Y}_t = \frac{1}{T}\left[\sum_{i=1}^N \hat{w}_t^i\left(h_{t+1}^T\left(\xi_{t+1}^i\right) - K h_{t+1}^T\left(\hat{\xi}_t^i\right)\right)\right],$$

*where we have introduced shorthand $h_t^T = \sum_{s=t}^T K^{s-t}\left(f - \mu\left(f\right)\right)$.*

PROOF. Use Lemma 4.1 and simplify terms. $\square$

Using the martingale in Lemma 4.3 we can prove the following lower bound on WE's asymptotic variance:

PROPOSITION 4.2. *Consider a WE scheme with a kernel $K$ that is geometrically ergodic and $V$-uniformly ergodic, with $\left\|f^2/V\right\| < \infty$. The variance of WE estimates for $\mu\left(f\right)$ is bounded from below by*

$$(4.36) \qquad \liminf_{T \to \infty} T \operatorname{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^N w_t^i f\left(\xi_t^i\right)\right] \geq \frac{\mu\left(v_f\right)^2}{N}.$$

PROOF. First, the $V$-uniform ergodicity and $\sqrt{V}$-uniform ergodicity of $K$ guarantee $\left\|v_f^2/V\right\| < \infty$. Hence, the right-hand side is finite and we can consider without loss of generality a subsequence of $T$ values for which $\operatorname{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^N w_t^i f\left(\xi_t^i\right)\right] < \infty$. Then a

martingale variance decomposition using Lemma 4.3 guarantees

$$(4.37) \qquad \mathrm{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N}w_t^i f\left(\xi_t^i\right)\right] \geq \sum_{t=0}^{T-2}\mathrm{E}\left|Y_{t+1}-\hat{Y}_t\right|^2 = \sum_{t=0}^{T-2}\mathrm{E}\left[\mathrm{Var}\left[Y_{t+1}\middle|\hat{\mathcal{F}}_t\right]\right].$$

Applying Jensen's inequality and setting $v_t^T(x) = \mathrm{Var}_{K(x,\cdot)}\left[h_{t+1}^T\right]^{1/2}$, we calculate

$$(4.38) \qquad T^2\,\mathrm{E}\left[\mathrm{Var}\left[Y_{t+1}\middle|\hat{\mathcal{F}}_t\right]\right] = \mathrm{E}\left[\mathrm{Var}\left[\sum_{i=1}^{N}w_{t+1}^i v_{t+1}^T\left(\xi_{t+1}^i\right)\middle|\hat{\mathcal{F}}_t\right]\right]$$

$$(4.39) \qquad = \mathrm{E}\left[\sum_{i=1}^{N}\left|\hat{w}_t^i v_t^T\left(\hat{\xi}_t^i\right)\right|^2\right]$$

$$(4.40) \qquad \geq \frac{1}{N}\mathrm{E}\left|\sum_{i=1}^{N}\hat{w}_t^i v_t^T\left(\hat{\xi}_t^i\right)\right|^2$$

$$(4.41) \qquad \geq \frac{\mathrm{E}\left[\sum_{i=1}^{N}\hat{w}_t^i v_t^T\left(\hat{\xi}_t^i\right)\right]^2}{N}$$

$$(4.42) \qquad = \frac{\mu_0\left(K^t v_t^T\right)^2}{N}.$$

In summary, we find

$$(4.43) \qquad T\,\mathrm{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N}w_t^i f\left(\xi_t^i\right)\right] \geq \frac{1}{TN}\sum_{t=0}^{T-1}\mu_0\left(K^t v_t^T\right)^2$$

$$(4.44) \qquad = \frac{1}{N}\int_0^1 \mu_0\left(K^{\lfloor sT\rfloor}v_{\lfloor sT\rfloor}^T\right)^2 ds.$$

For any $0 < s \leq 1$, we observe that $\left\|\mu_0 K^{\lfloor sT\rfloor}-\mu\right\| \overset{T\to\infty}{\to} 0$ and also $v_{\lfloor sT\rfloor}^T \overset{T\to\infty}{\to} v_f$ pointwise on the set $\{V < \infty\}$. Hence, by a useful generalization of Fatou's lemma (see [34, sec. 11.4]),

$$(4.45) \qquad \liminf_{T\to\infty}\mu_0\left(K^{\lfloor sT\rfloor}v_{\lfloor sT\rfloor}^T\right) \geq \mu\left(v_f\right).$$

We are able to conclude

$$(4.46) \qquad \liminf_{T\to\infty}T\,\mathrm{Var}\left[\frac{1}{T}\sum_{t=0}^{T-1}\sum_{i=1}^{N}w_t^i f\left(\xi_t^i\right)\right] \geq \frac{\mu\left(v_f\right)^2}{N}.$$

$\square$

Throughout our variance analysis, we have made the minimal assumptions that are needed to prove our results. In Proposition 4.2, we needed the assumption $\left\|f^2/V\right\| < \infty$ to ensure that the variance function $v_f^2$ is well-defined on a set of full $\mu$ measure and $\mu\left(v_f\right) < \infty$. Moving forward, in order to prove Lemma 3.1, we also need the assumption $\mu_0\left(V\right) < \infty$. This condition rules out a degenerate situation where the initial particles are drawn so far out of equilibrium that there is a lingering effect on the first and second moments of WE estimates — the same assumption would also be needed to bound the variance of direct MCMC estimates as well.

As we demonstrate below, our minimal assumptions are enough to verify Lemma 3.1, which gives a precise expression for WE's asymptotic variance.

PROOF OF LEMMA 3.1. We manipulate the martingale differences in Lemma 4.3 to find

$$(4.47) \qquad T^2 \operatorname{E} |Y_{t+1} - Y_t|^2 = T^2 \operatorname{E} \left[ \operatorname{Var} \left[ Y_{t+1} | \mathcal{F}_t \right] \right]$$

$$(4.48) \qquad = \operatorname{E} \left[ \operatorname{Var} \left[ \sum_{i=1}^{N} w_{t+1}^i h_{t+1}^T \left( \xi_{t+1}^i \right) \middle| \mathcal{F}_t \right] \right]$$

$$(4.49) \qquad = \operatorname{E} \left[ \sum_{i=1}^{N} \left| w_{t+1}^i \right|^2 \operatorname{Var} \left[ h_{t+1}^T \left( \xi_{t+1}^i \right) \middle| \mathcal{F}_t \right] \right]$$

$$(4.50) \qquad = \operatorname{E} \left[ \sum_u \frac{w_t(u)^2}{N_t(u)} \operatorname{Var}_{\eta_t^u K} \left[ h_{t+1}^T \right] \right],$$

where we have used the definition of binned multinomial resampling and we have set $\eta_t^u = \frac{1}{w_t(u)} \sum_{i \in u} w_t^i \delta \left( \xi_t^i \right)$. Hence,

$$(4.51)$$
$$\operatorname{Var} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f \left( \xi_t^i \right) \right] = \frac{\operatorname{Var}_{\mu_0} \left[ h_0^T \right]}{N T^2} + \frac{1}{T^2} \sum_{t=0}^{T-2} \operatorname{E} \left[ \sum_u \frac{w_t(u)^2}{N_t(u)} \operatorname{Var}_{\eta_t^u K} \left[ h_{t+1}^T \right] \right].$$

In this decomposition, the $\sqrt{V}$-uniform ergodicity of $K$ and the condition $\mu_0(V) < \infty$ guarantee the first term is asymptotically $\mathcal{O}\left(T^{-2}\right)$. To analyze the second term, we first calculate

$$(4.52) \qquad \left| \operatorname{E} \left[ \sum_u \frac{w_t(u)^2}{N_t(u)} \left( \operatorname{Var}_{\eta_t^u K} \left[ h_{t+1}^T \right] - \operatorname{Var}_{\eta_t^u K} \left[ h_f \right] \right) \right] \right|$$

$$(4.53) \qquad \leq \operatorname{E} \left[ \sum_u w_t(u) \left| \operatorname{Var}_{\eta_t^u K} \left[ h_{t+1}^T \right] - \operatorname{Var}_{\eta_t^u K} \left[ h_f \right] \right| \right]$$

$$(4.54) \qquad \leq \operatorname{E} \left[ \sum_u w_t(u) \operatorname{Var}_{\eta_t^u K} \left[ h_{t+1}^T + h_f \right] \right]^{1/2} \operatorname{E} \left[ \sum_u w_t(u) \operatorname{Var}_{\eta_t^u K} \left[ h_{t+1}^T - h_f \right] \right]^{1/2}$$

$$(4.55) \qquad \leq \operatorname{E} \left[ \sum_u w_t(u) \eta_t^u K \left| h_{t+1}^T + h_f \right|^2 \right]^{1/2} \operatorname{E} \left[ \sum_u w_t(u) \eta_t^u K \left| h_{t+1}^T - h_f \right|^2 \right]^{1/2}$$

$$(4.56) \qquad = \left( \mu_0 K^{t+1} \left| h_{t+1}^T + h_f \right|^2 \right)^{1/2} \left( \mu_0 K^{t+1} \left| h_{t-1}^T - h_f \right|^2 \right)^{1/2}.$$

This leads to the bound

$$(4.57) \qquad \left| \frac{1}{T^2} \sum_{t=0}^{T-2} \operatorname{E} \left[ \sum_u \frac{w_t(u)^2}{N_t(u)} \operatorname{Var}_{\eta_t^u K} \left[ h_{t+1}^T \right] \right] - \frac{1}{T^2} \sum_{t=0}^{T-2} \operatorname{E} \left[ \sum_u \frac{w_t(u)^2}{N_t(u)} \operatorname{Var}_{\eta_t^u K} \left[ h_f \right] \right] \right|$$

$$(4.58) \qquad \leq \frac{1}{T^2} \sum_{t=0}^{T-2} \left( \mu_0 K^{t+1} \left| h_{t+1}^T + h_f \right|^2 \right)^{1/2} \left( \mu_0 K^{t+1} \left| h_{t+1}^T - h_f \right|^2 \right)^{1/2}.$$

The $\sqrt{V}$-uniform ergodicity of $K$, the $V$-uniform ergodicity of $K$, and the condition $\mu_0(V) < \infty$ guarantee that the last quantity is $\mathcal{O}\left(T^{-2}\right)$ as $T \to \infty$, confirming the result. $\square$

As the last step in our technical analysis, we use Lemma 3.1 to construct a WE scheme that nearly achieves the optimal variance bound.

PROPOSITION 4.3. *Consider a WE scheme with a kernel $K$ that is geometrically ergodic and $V$-uniformly ergodic, with $\left\| f^2/V \right\| < \infty$. Then, for any $\epsilon > 0$, there is a WE scheme that satisfies*

$$(4.59) \qquad \limsup_{T \to \infty} T \operatorname{Var} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right) \right] \leq (1+\epsilon) \frac{\mu\left(v_f\right)^2}{N}$$

*if the number of particles $N$ is sufficiently large.*

PROOF. In the case $\mu\left(v_f\right) = 0$, we must also have $\mu\left(v_f^2\right) = 0$, and direct MCMC sampling is sufficient to achieve the asymptotic variance upper bound. Next, we consider the case $\mu\left(v_f\right) > 0$. We assume initial particles are drawn from a distribution satisfying $\mu_0\left(V\right) < \infty$. We define bins based on spatial sets

$$(4.60) \qquad u_{i,j} = \left\{ x \in X : i - \frac{1}{2} < \frac{Kh_f\left(x\right)}{\delta} \leq i + \frac{1}{2}, \quad j - \frac{1}{2} < \frac{v_f\left(x\right)}{\delta} \leq j + \frac{1}{2} \right\},$$

$$(4.61) \qquad u_\infty = X \setminus \left(u_{i,j}\right)_{-J \leq i,j \leq J}$$

where $\delta$ and $J$ are parameters to be tuned. Here, in a slight abuse of notation, we are using $u_{i,j}$ to refer both to a spatial set and to the indices of the particles in that set. We set bin allocations $N_t\left(u\right)$ to satisfy

$$(4.62) \qquad \frac{N_t\left(u\right)}{N} \geq \max\left\{ \delta w_t\left(u\right), \left(1 - 2\delta\right) \frac{w_t\left(u\right) \eta_t^u\left(v_f\right)}{\sum_u w_t\left(u\right) \eta_u^u\left(v_f\right)} \right\},$$

which is always possible when the number of particles $N$ is sufficiently large.

Having introduced an explicit WE scheme, we bound its asymptotic variance using Lemma 3.1. We perform the following three-step variance calculation:

*Step 1.* We bound the intrabin variance in the $Kh_f$ and $v_f$ coordinates using

$$(4.63) \qquad \frac{1}{T} \sum_{t=0}^{T-2} \operatorname{E}\left[ \sum_u \frac{w_t\left(u\right)^2}{N_t\left(u\right)} \left(\operatorname{Var}_{\eta_t^u}\left[Kh_f\right] + \operatorname{Var}_{\eta_t^u}\left[v_f\right]\right) \right]$$

$$(4.64) \qquad \leq \frac{\delta}{2N} + \frac{1}{\delta TN} \sum_{t=0}^{T-2} \operatorname{E}\left[ w_t\left(u_\infty\right) \left(\operatorname{Var}_{\eta_t^{u_\infty}}\left[Kh_f\right] + \operatorname{Var}_{\eta_t^{u_\infty}}\left[v_f\right]\right) \right]$$

$$(4.65) \qquad \leq \frac{\delta}{2N} + \frac{1}{\delta TN} \sum_{t=0}^{T-2} \operatorname{E}\left[ w_t\left(u_\infty\right) \eta_t^{u_\infty} \left(\left|Kh_f\right|^2 + \left|v_f\right|^2\right) \right]$$

$$(4.66) \qquad \leq \frac{\delta}{2N} + \frac{\mu\left(\mathbb{1}_{u_\infty}\left(\left(Kh_f\right)^2 + v_f^2\right)\right)}{\delta N}$$

*Step 2.* We bound the remaining asymptotic variance term by using

$$(4.67) \qquad \frac{1}{T} \sum_{t=0}^{T-2} \operatorname{E}\left[ \sum_u \frac{\left|w_t\left(u\right) \eta_t^u\left(v_f\right)\right|^2}{N_t\left(u\right)} \right]$$

$$(4.68) \qquad \leq \frac{1}{\left(1 - 2\delta\right) NT} \sum_{t=0}^{T-2} \operatorname{E}\left| \sum_{i=1}^{N} w_t^i v_f\left(\xi_t^i\right) \right|^2$$

$$(4.69) \qquad = \frac{1}{\left(1 - 2\delta\right) NT} \sum_{t=0}^{T-2} \left( \mu\left(v_f\right)^2 + \operatorname{Var}\left[ \sum_{i=1}^{N} w_t^i v_f\left(\xi_t^i\right) \right] \right).$$

*Step 3.* To bound a quantity $\mathrm{Var}\left[\sum_{i=1}^N w_T^i v_f\left(\xi_T^i\right)\right]$, we consider the martingale $M_t$ that was introduced in Lemma 4.1. Using the function $v_f$ in place of $f$, Lemma 4.1 yields:

$$(4.70) \qquad \mathrm{Var}\left[\sum_{i=1}^N w_T^i v_f\left(\xi_T^i\right)\right]$$

$$(4.71) \qquad = \mathrm{Var}\left[\frac{1}{N}\sum_{i=1}^N K^T v_f\left(\xi_0^i\right)\right] + \sum_{t=0}^{T-1}\mathrm{E}\left[\mathrm{Var}\left[\sum_{i=1}^N w_{t+1}^i K^{T-t-1} v_f\left(\xi_t^i\right)\Big|\mathcal{F}_t\right]\right]$$

$$(4.72) \qquad = \frac{1}{N}\mathrm{Var}_\mu\left[K^T v_f\right] + \sum_{t=0}^{T-1}\mathrm{E}\left[\sum_u \frac{w_t(u)^2}{N_t(u)}\mathrm{Var}_{\eta_t^u K}\left[K^{T-t-1} v_f\right]\right]$$

$$(4.73) \qquad \leq \frac{1}{\delta N}\mathrm{Var}_\mu\left[K^T v_f\right] + \frac{1}{\delta N}\sum_{t=0}^{T-1}\mathrm{E}\left[\sum_{i=1}^N w_t^i \mathrm{Var}_{K(\xi_t, \cdot)}\left[K^{T-t-1} v_f\right]\right]$$

$$(4.74) \qquad = \frac{1}{\delta N}\sum_{t=0}^T \mathrm{Var}_\mu\left[K^t v_f\right]$$

$$(4.75) \qquad \leq \frac{1}{\delta N}\sum_{t=0}^\infty \mathrm{Var}_\mu\left[K^t v_f\right].$$

We confirm this last term is finite, because $\left\|v_f/\sqrt{V}\right\| < \infty$ and $K$ is $\sqrt{V}$-uniformly ergodic.

In summary, steps 1-3 reveal that

$$(4.76) \qquad \frac{1}{T}\sum_{t=0}^{T-2}\mathrm{E}\left[\sum_u \frac{w_t(u)^2}{N_t(u)}\left(\mathrm{Var}_{\eta_t^u}[Kh] + \mathrm{Var}_{\eta_t^u}[v_f] + \eta_t^u(v_f)^2\right)\right]$$

$$(4.77) \qquad \leq \frac{\delta}{2N} + \frac{\mu\left(\mathbb{1}_{u_\infty}\left((Kh_f)^2 + v_f^2\right)\right)}{\delta N} + \frac{\mu(v_f)^2}{(1-2\delta)N} + \sum_{t=0}^\infty \frac{\mathrm{Var}_\mu\left[K^t v_f\right]}{(\delta - 2\delta^2)N^2}$$

By taking $\delta$ appropriately small and then taking $J$ and $N$ appropriately large, we can make this last quantity less than $(1+\epsilon)\mu(v_f)^2/N$, thereby completing the proof. $\qquad\square$

**5. Numerical experiments.** In this section, we apply WE to compute rare event probabilities in three example problems. These numerical experiments validate our formulas for WE's optimal variance while also demonstrating the major potential for efficiency gains by using WE instead of MCMC.

**5.1. *Geometric tail probabilities.*** In the first example, our goal is estimating tail probabilities $\mu[a, \infty)$ for the geometric distribution

$$(5.1) \qquad \mu(x) = 2^{-x-1}, \quad x \in \mathbb{Z}^+ = \{0, 1, \ldots\}.$$

To sample from $\mu$, we use a Markov chain with transition probabilities

$$(5.2) \qquad P(x, x+1) = P(x, 0) = \frac{1}{2}.$$

When $a$ is large, it would be very costly to estimate tail probabilities $\mu[a, \infty) = 2^{-a}$ by direct MCMC sampling. However, we show that WE can make these calculations more tractable.

5.1.1. *WE implementation.* In our numerical experiments, we use WE to estimate the tail probability $\mu[a, \infty) = 2^{-a}$ for $a = 25$. We draw initial particles from $\mu$, and we sample for $T = 1000$ time steps. Following the optimization strategy discussed in Section 3.2, we sort the particles into spatial bins based on the sets

(5.3) $$u_i = \{i\}, \quad 0 \leq i \leq 23, \quad u_{24} = [24, \infty),$$

which are the exact level sets of $Kh_f$. Then, we allocate children particles to each bin according to the rule

(5.4) $$\frac{N_t(u)}{N} \approx \frac{w_t(u)\eta_t^u(v_f)}{\sum_{u'} w_t(u')\eta_t^{u'}(v_f)}.$$

5.1.2. *WE results.* In Figure 5 below, we present WE's relative variance constant

(5.5) $$\text{Relative Variance Constant} = \frac{NT}{\mu(f)^2} \text{Var}\left[\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{N} w_t^i f(\xi_t^i)\right],$$

calculated over $10^6$ independent trials for the function $f(x) = \mathbb{1}\{x \geq 25\}$. Additionally, we present theoretical relative variance constants for MCMC and WE, calculated using the asymptotic theory developed in Section 3. With just $N = \mathcal{O}(a)$ particles, we find that WE very nearly achieves the theoretical optimal variance, thereby improving MCMC's variance by more than five orders of magnitude.
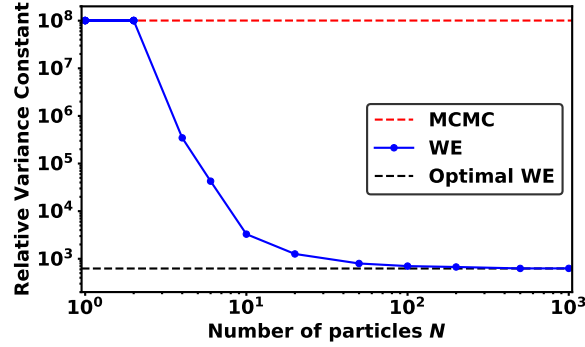


FIG 5. *Application of WE to the geometric tails problem.*

5.2. *Gaussian tail probabilities.* In our second example, we use WE to estimate tail probabilities $\mu[a, \infty)$ for the Gaussian distribution $\mu = \mathcal{N}(0, 1)$. To sample from $\mu$, we use the first-order autoregressive process

(5.6) $$X_{k+1} = e^{-\Delta t}X_k + \sqrt{1 - e^{-2\Delta t}}\eta_{k+1}, \quad \eta_{k+1} \sim N(0, 1).$$

5.2.1. *WE implementation.* In our numerical tests, we apply WE to estimate the tail probabilities $\mu[3, \infty) = 1.35 \times 10^{-3}$ and $\mu[4, \infty) = 3.17 \times 10^{-5}$. We start all the particles at $x = 0$, and then we simulate forward for $n_T = T/\Delta t$ time steps, where $T = 10^4$ and $\Delta t = 0.01$. At each splitting step, we sort the particles into bins based on the intervals $(x_i, x_{i+1}]$, where

(5.7) $$-\infty = x_0 < x_1 < \cdots < x_{\max-1} < x_{\max} = \infty.$$

We optimize the mesh points $x_2, \ldots, x_{\max-2}$ to ensure that intervals $(x_i, x_{i+1}]_{1 \leq i \leq \max-2}$ are approximate level sets of $h_f$. We use the `WeightedEnsemble.jl` package [3] for our numerical implementation and describe additional implementation details in Appendix A.

5.2.2. *WE error bars.* In this example, we consider two data-driven strategies for estimating the variance of WE estimates. As a first strategy, we run WE for $100$ independent trials and apply a bootstrap approach for estimating the variance [5, 14]. In this bootstrap approach, we generate $M = 10^4$ bootstrap samples of size $100$ by randomly subsampling from the independent WE estimates. Then, for each bootstrap sample, we compute the empirical variance. By aggregating together the $M = 10^4$ variance estimates, we obtain a point estimate and robust confidence intervals for WE's variance.

As a second strategy for variance estimation, we apply the following variance estimate to each one of the independent WE runs:

$$(5.8) \qquad \mathrm{Var}\left[ \frac{1}{n_T} \sum_{t=0}^{n_T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right) \right]$$

$$(5.9) \qquad \approx \frac{1}{n_T^2} \sum_{|t-s| \leq L} \left( \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right) - \hat{\mu}(f) \right) \left( \sum_{i=1}^{N} w_s^i f\left(\xi_s^i\right) - \hat{\mu}(f) \right).$$

In this formula,

$$(5.10) \qquad \hat{\mu}(f) = \frac{1}{n_T} \sum_{t=0}^{n_T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)$$

is the empirical estimate of $\mu(f)$, while $L \geq 0$ is a truncation threshold, chosen so that correlations between $\sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right)$ and $\sum_{i=1}^{N} w_s^i f\left(\xi_s^i\right)$ are negligible for any time lag $|s - t|$ exceeding $L$.

The variance estimator (5.9) is potentially very useful, since it provide error bars for WE estimates even after a single run of the algorithm. Indeed, (5.9) is already the standard variance estimator in MCMC, and among MCMC practitioners it is known as the integrated autocorrelation time (IAT) estimator [38]. When the IAT estimator is applied to WE results, the full convergence properties have not yet been rigorously guaranteed. However, we observe that the estimator has asymptotic bias that vanishes exponentially fast as we increase the truncation threshold $L$. Moreover, in our experiments, we find good agreement between variance estimates using the IAT estimator and those obtained using the bootstrap. Our results provide empirical evidence that, at least for some problems to which the WE is applied, the IAT estimator is a useful tool.

5.2.3. *WE results.* In Figure 6, we present our estimates of the relative variance constant

$$(5.11) \qquad \text{Relative Variance Constant} = \frac{NT}{\mu(f)^2} \mathrm{Var}\left[ \frac{1}{n_T} \sum_{t=0}^{n_T-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right) \right],$$

where $f(x) = \mathbb{1}\{x \geq 3\}$ in the first scenario. and $f(x) = \mathbb{1}\{x \geq 4\}$ in the second scenario. We compare our experimental estimates against asymptotic formulas for the relative variance constant that are valid in the simultaneous limit as $T \to \infty$ and $\Delta t \to 0$. A full derivation of these formulas appears in the appendix.

For $a = 3$ and sufficiently large $N$, WE nearly attains the optimal variance, improving MCMC's variance by over an order of magnitude. For $a = 4$, WE's variance is somewhat further from the optimal variance, yet WE still achieves over two orders of magnitude improvement over direct MCMC sampling.
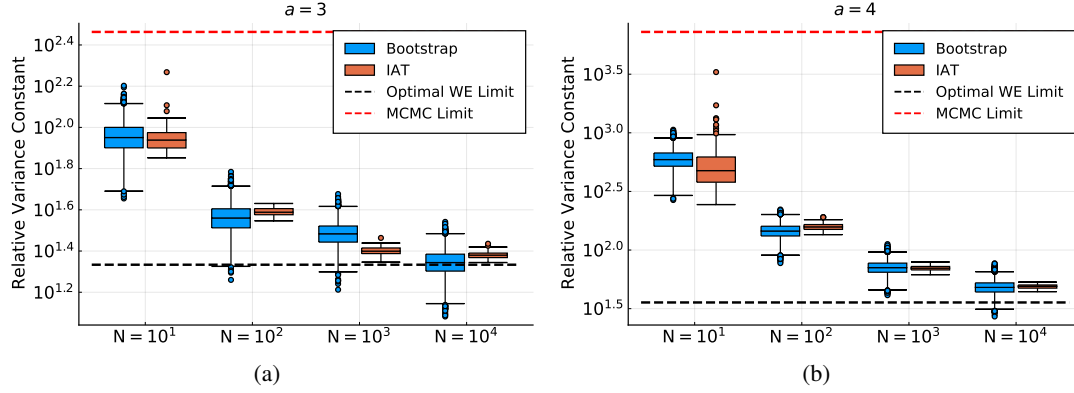
FIG 6. *Application of WE to the Gaussian tails problem.*

5.3. *Ising tail probabilities.* In our third and final example, we use WE to calculate the probability of extreme magnetizations for the Ising model on an $L \times L$ lattice with periodic boundary conditions. The Ising model has long been the subject of study in the statistical physics community as a model of ferromagnetism and as simple system exhibiting phase changes [6, 19, 37]. The energy associated with the model is

$$(5.12) \qquad H(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{i \sim j} \boldsymbol{\sigma}_i \boldsymbol{\sigma}_j, \quad \boldsymbol{\sigma}_i \in \{+1, -1\},$$

where $i \sim j$ denotes that $i$ and $j$ are neighboring lattice points. The associated Boltzmann distribution is

$$(5.13) \qquad \mu(\boldsymbol{\sigma}) = \frac{\exp(-\beta H(\boldsymbol{\sigma}))}{Z}, \quad Z = \sum_{\boldsymbol{\sigma}'} \exp\left(-\beta H\left(\boldsymbol{\sigma}'\right)\right).$$

When $\beta > \beta_c$ (the "low-temperature" regime), the system tends to self-organize with the majority of spins all either $+1$ or all $-1$. On the other hand, when $\beta < \beta_c$ (the "high-temperature" regime), self-organization is less likely, and a mixture of $+1$s and $-1$s becomes more likely.

Our numerical tests address the following questions:

- What is the probability that the mean magnetization, $m(\boldsymbol{\sigma}) = L^{-2} \sum_i \boldsymbol{\sigma}_i$, is in $(-0.1, 0.1)$ in the low-temperature regime? In other words, what is the likelihood of seeing the system in a highly disordered state, despite being at low temperature?
- What is the probability that the mean magnetization satisfies $|m| > 0.9$ in the high-temperature regime? Here, we are considering the likelihood of seeing the system in a highly ordered state, despite being at high temperature.

5.3.1. *WE implementation.* In our experiments, we implement WE on a $10 \times 10$ lattice. In the low-temperature regime, we start 100 particles from an initial state of all $-1$s. In the high-temperature regime, we start 100 particles from an initial state randomly selected from the uniform distribution on spins. In both regimes, we evolve the particles forward by selecting one of the $L^2$ spins uniformly and proposing a flip from $\boldsymbol{\sigma}_i$ to $-\boldsymbol{\sigma}_i$. We accept this proposed change with probability

$$(5.14) \qquad \min\left\{1, \exp\left(-\beta \boldsymbol{\sigma}_i \sum\nolimits_{j \sim i} \boldsymbol{\sigma}_j\right)\right\},$$

and otherwise leave the system unchanged. We perform ten such updates in each forward evolution step. Then, in each splitting step, we sort particles into bins based on mean magnetization and apply splitting and killing. We describe additional details in Appendix B.

5.3.2. *WE results.*   In Figures 7 and 8, we report the mean and standard deviation of the running averages

$$(5.15) \qquad \frac{1}{t} \sum_{s=0}^{t-1} \sum_{i=1}^{N} w_t^i f\left(\xi_t^i\right), \quad t = 0, 1, \dots T-1,$$

computed over 100 independent trials for the functions $f(\boldsymbol{\sigma}) = \mathbb{1}\{|m(\boldsymbol{\sigma})| > 0.9\}$ and $f(\boldsymbol{\sigma}) = \mathbb{1}\{|m(\boldsymbol{\sigma})| < 0.1\}$. We also report the relative variance constants based on the running averages.



FIG 7. *Application of WE to the Ising model at a low temperature ($\beta = 0.6$).*
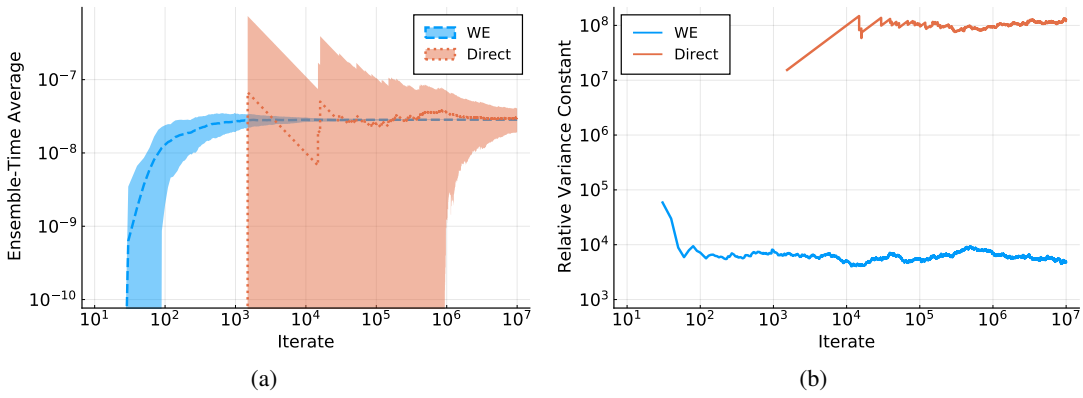


FIG 8. *Application of WE to the Ising model at a high temperature ($\beta = 0.25$).*

Not only do we find that WE is more computationally efficient than MCMC, but our results also show that WE is more efficient than sampling from the Ising model by using an independence sampler. An independence sampler would lead to a relative variance constant of $p^{-1} - 1$ when estimating a rare probability $p$. Yet Figures 7 and 8 show that WE improves this variance constant by several orders of magnitude, providing especially large improvements in the high-temperature regime. In conclusion, we obtain a remarkable result: WE transforms the time correlations in the dynamics, which would normally be an impediment to efficient sampling [38], into a major asset that enables significant variance reduction.

**6. Conclusion.** In this work, we presented splitting as an approach for reducing MCMC's variance when estimating rare event probabilities. Traditionally, splitting is viewed as separate from MCMC in the mathematical literature. However, here we showed that splitting can be beneficially combined with MCMC when appropriate stability conditions are satisfied. We contributed the following results:

1. We showed that splitting schemes can degenerate over long timescales due to shrinking weights. Moreover, we proved that the only way to avoid shrinking weights is by using weighted ensemble (WE).
2. We presented an optimal variance bound for WE that demonstrates the method's maximal efficiency when a large number of particles are available.
3. We explored numerical examples where WE reduces MCMC's variance by multiple orders of magnitude.

As our numerical examples make clear, there remain significant open questions for investigation. First, it would be desirable to estimate the variance of WE estimates from a single long trajectory of WE data. Yet it remains to be determined whether the integrated autocorrelation time (IAT) estimator provides convergent estimates of WE's variance. Second, it is clear from our examples that WE requires a large number of particles in order to attain peak efficiency. The precise scaling of the variance with the number of particles is an open area of investigation.

In light of these open questions, we regard our present work not as the final answer regarding WE's properties but rather as an essential step toward uncovering the method's mathematical foundations. Here, we have demonstrated WE's importance as a practical computational tool and its interest as a mathematical system where interactions perturb the behavior of ergodic Markov chains. We have shown that despite the apparent complexity of WE's dynamics, the mean and variance of WE's estimates can be precisely bounded, yielding insights into the method's efficiency. In summary, we have established the unique role of WE as a splitting method that reduces MCMC variance and constructed a rigorous framework that will aid in the method's future development.

## APPENDIX A: DETAILS OF OU COMPUTATIONS

To calculate WE's optimal variance, we use asymptotic approximations that are valid in the simultaneous limit as $T \to \infty$ and $\Delta t \to 0$. We observe that the process (5.6) is the $\Delta t$-skeleton of the continuous-time Ornstein-Uhlenbeck (OU) process

$$(A.1) \qquad d\bar{X}_t = -\bar{X}_t \, dt + \sqrt{2} \, d\bar{W}_t \, .$$

Therefore, when $\Delta t \ll 1$, we can approximate the conditional expectation function $h_f$ using

$$(A.2) \qquad \bar{h}_f = \frac{1}{\Delta t} \, \mathrm{E}_x \left[ \int_0^\infty \mathbb{1} \left\{ \bar{X}_t \geq a \right\} - \mu \left[ a, \infty \right) dt \right] .$$

Likewise, we can approximate the variance function $v_f^2 = Kh_f^2 - (Kh_f)^2$ using

$$(A.3) \qquad \bar{v}_f^2(x) = \Delta t \lim_{t \to 0+} \frac{1}{t} \, \mathrm{E}_x \left| \bar{h}_f(X_t) - \bar{h}_f(X_0) \right|^2 .$$

The approximation as $\Delta t \ll 0$ leads to useful simplifications, since $\bar{v}_f^2$ is determined by the quadratic variation [28] of the process $\bar{h}_f(X_t)$; hence,

$$(A.4) \qquad \bar{v}_f^2(x) = 2\Delta t \left| \frac{d\bar{h}_f(x)}{dx} \right|^2$$

To calculate the conditional expectation function $\overline{h}_f$ and the variance function $\overline{v}_f^2$ explicitly using Mathematica, we first observe that

$$(A.5) \qquad \Delta t \overline{h}_f = \int_0^\infty \left( P_x \left\{ \overline{X}_t \geq a \right\} - \mu\left[a, \infty\right) \right) dt$$

solves the Poison equation

$$(A.6) \qquad -\mathcal{L}\left(\Delta t \overline{h}_f\right) = \mathbb{1}\left\{x \geq a\right\} - \mu\left[a, \infty\right)$$

involving the infinitesimal generator of the OU process $\mathcal{L}g = -xg' + g''$. Hence, the approximate variance function $\overline{v}_f^2 = 2\Delta t \left|\overline{h}_f'\right|^2$ solves the first-order ODE

$$(A.7) \qquad x\overline{v}_f - \overline{v}_f' = \sqrt{\frac{2}{\Delta t}}\left(\mathbb{1}\left\{x \geq a\right\} - \mu\left[a, \infty\right)\right).$$

Solving the ODE gives

$$(A.8) \qquad \overline{v}_f(x) = \sqrt{\frac{2}{\Delta t}} \frac{\min\left\{\Phi(x), \Phi(a)\right\} - \Phi(x)\Phi(a)}{\phi(x)},$$

where

$$(A.9) \qquad \phi(x) = \frac{\exp\left(-x^2/2\right)}{\sqrt{2\pi}}, \quad \text{and} \quad \Phi(x) = \int_{-\infty}^x \phi(y)\, dy$$

are the probability density function and cumulative distribution function for a Gaussian distribution. Using formula (A.8), we conclude that the MCMC variance and the optimal WE variance can be approximated as follows.

$$(A.10) \qquad \text{MCMC variance:} \qquad \frac{\mu\left(\overline{v}_f^2\right)}{NT/\Delta t} = \frac{4\exp\left(-a^2/2\right)}{\sqrt{2\pi}a^3 NT}\left(1 + \mathcal{O}\left(a^{-2}\right)\right).$$

$$(A.11) \qquad \text{Optimal WE variance:} \qquad \frac{\mu\left(\overline{v}_f\right)^2}{NT/\Delta t} = \frac{\exp\left(-a^2\right)}{\pi NT}.$$

Thus, we find that the optimal improvement factor of WE over MCMC increases exponentially fast as $a \to \infty$. Lastly, using Mathematica we integrate (A.8) to obtain a closed-form expression for $\overline{h}_f$ involving confluent hypergeometric functions of the first kind.

In our implementation of WE, we define bins using a mesh

$$(A.12) \qquad -\infty = x_0 < x_1 < \cdots < x_{\max-1} < x_{\max}. = \infty.$$

The endpoints of the mesh are set to $x_1 = -2$ and $x_{\max} = 3.5$ in the case $a = 3$, and $x_1 = -2$ and $x_{\max} = 5$ in the case $a = 4$. The interior mesh points $x_2, x_3, \ldots, x_{\max-1}$ are chosen to constrain the variation of $\Delta t \overline{h}_f$ over each of the intervals $(x_i, x_{i+1}]_{1 \leq i \leq \max-2}$. The variation per interval is set to $10^{-3}$ in the case $a = 3$ and $10^{-4}$ in the case $a = 4$.

Lastly, during the WE run, we allocate children particles to each bin according to the rule

$$(A.13) \qquad \frac{N_t(u)}{N} \approx \frac{w_t(u)\,\eta_t^u\left(\overline{v}_f\right)}{\sum_{u'} w_t(u')\,\eta_t^{u'}\left(\overline{v}_f\right)},$$

as described in [4]. We use systematic resampling to select particles within the bins.

## APPENDIX B: DETAILS OF ISING COMPUTATIONS

We set the bins to be Voronoi cells in the magnetization coordinate $m$ with centers $-1, -0.9, \ldots, 0.9, 1$. We allocate children particles to each bin according to the rule

$$
(\text{B.1}) \qquad \frac{N_t(u)}{N} \approx \frac{w_t(u)\,\eta_t^u(\overline{v}_f)}{\sum_{u'} w_t(u')\,\eta_t^{u'}(\overline{v}_f)},
$$

where $\overline{v}_f$ is an approximation to $v_f$ built on a coarse model of the dynamics.

To obtain $\overline{v}_f$, we follow the microbin approach developed in [1, 4]. We first use short, independent simulations to obtain a transition matrix $\overline{K}$ for the coordinate $m$. Specifically, by sampling from the uniform distribution with fixed magnetization $m$, we obtain $10^4$ initial data points in each magnetization state

$$
(\text{B.2}) \qquad m = -1, -1 + 2L^{-2}, \ldots 1 - 2L^{-2}, 1.
$$

Then, we run the dynamics forward for one evolution step to estimate the entries

$$
(\text{B.3}) \qquad \overline{K}_{ij} = \sum_{m(\boldsymbol{\sigma})=i} \mu(\boldsymbol{\sigma})\,K(\boldsymbol{\sigma}, \mathbb{1}\{m = m_j\}).
$$

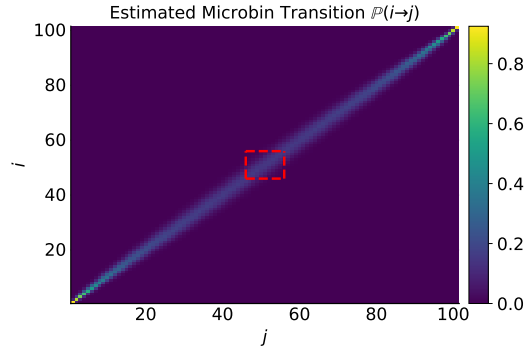We show the estimated $\overline{K}$ matrix in Figure 9 below.



FIG 9. *Microbin transition matrix for the low-temperature Ising model. The red square indicates the low-magnetization state $|m| < 0.1$ whose probability we seek to estimate.*

Having obtained $\overline{K}$, the microbin transition matrix, we next compute the microbin invariant measure $\overline{\mu}^T = \overline{\mu}^T \overline{K}$. Lastly, we solve the Poisson equation

$$
(\text{B.4}) \qquad (I - \overline{K})\,\overline{h}_f = f - \overline{u}(f)
$$

to approximate the conditional expectation function $\overline{h}_f$ and the variance function $\overline{v}_f(x)^2 = \mathrm{Var}_{\overline{K}(x,\cdot)}\left[\overline{h}_f\right]$.

## REFERENCES

[1] ARISTOFF, D. (2018). Analysis and optimization of weighted ensemble sampling. *ESAIM: Mathematical Modelling and Numerical Analysis* **52** 1219–1238.

[2] ARISTOFF, D. (2019). An ergodic theorem for weighted ensemble. *arXiv preprint arXiv:1906.00856*.

[3] ARISTOFF, D., JONES, F. G., WEBBER, R. J., SIMPSON, G. and ZUCKERMAN, D. M. (2020). WeightedEnsemble.jl. Julia package.

[4] ARISTOFF, D. and ZUCKERMAN, D. M. (2020). Optimizing weighted ensemble sampling of steady states. *Multiscale Modeling & Simulation* **18** 646–673.

[5] ASMUSSEN, S. and GLYNN, P. W. (2007). *Stochastic simulation: Algorithms and analysis* **57**. Springer Science & Business Media.

[6] BAXTER, R. J. (2016). *Exactly solved models in statistical mechanics*. Elsevier.

[7] BHATT, D., ZHANG, B. W. and ZUCKERMAN, D. M. (2010). Steady-state simulations using weighted ensemble path sampling. *The Journal of chemical physics* **133** 014110.

[8] CÉROU, F. and GUYADER, A. (2007). Adaptive multilevel splitting for rare event analysis. *Stochastic Analysis and Applications* **25** 417–443.

[9] CÉROU, F., GUYADER, A. and ROUSSET, M. (2019). Adaptive multilevel splitting: Historical perspective and recent results. *Chaos: An Interdisciplinary Journal of Nonlinear Science* **29** 043108.

[10] COPPERMAN, J. and ZUCKERMAN, D. M. (2020). Accelerated estimation of long-timescale kinetics from weighted ensemble simulation via non-Markovian "microbin" analysis. *Journal of Chemical Theory and Computation* **16** 6763–6775.

[11] COPPERMAN, J. T. and ZUCKERMAN, D. M. (2020). Accelerated estimation of long-timescale kinetics by combining weighted ensemble simulation with Markov model "microstates" using non-Markovian theory. *Biophysical Journal* **118** 180a.

[12] COSTAOUEC, R., FENG, H., IZAGUIRRE, J. and DARVE, E. (2013). Analysis of the accelerated weighted ensemble methodology. In *Conference Publications* **2013** 171. American Institute of Mathematical Sciences.

[13] DARVE, E. and RYU, E. (2012). Chapter 7. Computing reaction rates in bio-molecular systems using discrete macro-states. In *RSC Biomolecular Sciences* 138–206. Royal Society of Chemistry.

[14] DAVISON, A. C. and HINKLEY, D. V. (1997). *Bootstrap methods and their application* **1**. Cambridge university press.

[15] DEL MORAL, P. (2012). *Feynman-Kac formulae: Genealogical and interacting particle systems with applications*. Springer Science & Business Media.

[16] DEL MORAL, P., GARNIER, J. et al. (2005). Genealogical particle analysis of rare events. *The Annals of Applied Probability* **15** 2496–2534.

[17] DINNER, A. R., MATTINGLY, J. C., TEMPKIN, J. O., KOTEN, B. V. and WEARE, J. (2018). Trajectory stratification of stochastic dynamics. *SIAM Review* **60** 909–938.

[18] DOUC, R. and CAPPÉ, O. (2005). Comparison of resampling schemes for particle filtering. In *ISPA 2005. Proceedings of the 4th International Symposium on Image and Signal Processing and Analysis, 2005.* 64–69. IEEE.

[19] GALLAVOTTI, G. (2013). *Statistical mechanics: A short treatise*. Springer Science & Business Media.

[20] GEMAN, S. and GEMAN, D. (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on pattern analysis and machine intelligence* **6** 721–741.

[21] GLASSERMAN, P., HEIDELBERGER, P., SHAHABUDDIN, P. and ZAJIC, T. (1999). Multilevel splitting for estimating rare event probabilities. *Operations Research* **47** 585–600.

[22] GRASSBERGER, P. (1997). Pruned-enriched Rosenbluth method: Simulations of $\theta$ polymers of chain length up to 1000000. *Physical Review E* **56** 3682.

[23] GUBERNATIS, J., KAWASHIMA, N. and WERNER, P. (2016). *Quantum Monte Carlo methods*. Cambridge University Press.

[24] HUBER, G. A. and KIM, S. (1996). Weighted-ensemble Brownian dynamics simulations for protein association reactions. *Biophysical journal* **70** 97–110.

[25] HUSSAIN, S. and HAJI-AKBARI, A. (2020). Studying rare events using forward-flux sampling: Recent breakthroughs and future outlook. *The Journal of Chemical Physics* **152** 060901.

[26] JONES, G. L. et al. (2004). On the Markov chain central limit theorem. *Probability surveys* **1** 299–320.

[27] KAHN, H. and HARRIS, T. E. (1951). Estimation of particle transmission by random sampling. *National Bureau of Standards applied mathematics series* **12** 27–30.

[28] KALLENBERG, O. (2006). *Foundations of modern probability*. Springer Science & Business Media.

[29] LIU, J. S. (2008). *Monte Carlo strategies in scientific computing*. Springer Science & Business Media.

[30] METROPOLIS, N., ROSENBLUTH, A. W., ROSENBLUTH, M. N., TELLER, A. H. and TELLER, E. (1953). Equation of state calculations by fast computing machines. *The journal of chemical physics* **21** 1087–1092.

[31] MEYN, S. P. and TWEEDIE, R. L. (2012). *Markov chains and stochastic stability*. Springer Science & Business Media.

[32] PRATT, A., SUÁREZ, E., ZUCKERMAN, D. M. and CHONG, L. T. (2019). Extensive evaluation of weighted ensemble strategies for calculating rate constants and binding affinities of molecular association/dissociation processes. *bioRxiv* 671172.

[33] ROSENBLUTH, M. N. and ROSENBLUTH, A. W. (1955). Monte Carlo calculation of the average extension of molecular chains. *The Journal of Chemical Physics* **23** 356–359.

[34] ROYDEN, H. L. (1988). *Real analysis, 3rd edition*. Pearson Custom Publishing.

[35] RUBINO, G. and TUFFIN, B. (2009). *Rare event simulation using Monte Carlo methods*. John Wiley & Sons.

[36] RUBINSTEIN, R. Y. and KROESE, D. P. (2016). *Simulation and the Monte Carlo method*. John Wiley & Sons, Inc.

[37] RUELLE, D. (1999). *Statistical mechanics: Rigorous results*. World Scientific.

[38] SOKAL, A. (1997). Monte Carlo methods in statistical mechanics: Foundations and new algorithms. In *Functional integration* 131–192. Springer.

[39] TORRILLO, P. A., BOGETTI, A. T. and CHONG, L. T. (2020). A minimal, adaptive binning scheme for weighted ensemble simulations. *bioRxiv*.

[40] WEBBER, R. J. (2019). Unifying Sequential Monte Carlo with resampling matrices. *arXiv preprint arXiv:1903.12583*.

[41] ZHANG, B. W., JASNOW, D. and ZUCKERMAN, D. M. (2010). The "weighted ensemble" path sampling method is statistically exact for a broad class of stochastic processes and binning procedures. *The Journal of chemical physics* **132** 054107.

[42] ZUCKERMAN, D. M. and CHONG, L. T. (2017). Weighted ensemble simulation: Review of methodology, applications, and software. *Annual review of biophysics* **46** 43–57.