

A simple method for automated equilibration detection in molecular simulations

John D. Chodera^{1,*}

¹Computational Biology Program, Sloan Kettering Institute,
Memorial Sloan Kettering Cancer Center, New York, NY 10065

(Dated: December 15, 2015)

Molecular simulations intended to compute equilibrium properties are often initiated from configurations that are highly atypical of equilibrium samples, a practice which can generate a distinct initial transient in mechanical observables computed from the simulation trajectory. Traditional practice in simulation data analysis recommends this initial portion be discarded to *equilibration*, but no simple, general, and automated procedure for this process exists. Here, we suggest a conceptually simple automated procedure that does not make strict assumptions about the distribution of the observable of interest, in which the equilibration time is chosen to maximize the number of effectively uncorrelated samples in the production timespan used to compute equilibrium averages. We present a simple Python reference implementation of this procedure, and demonstrate its utility on typical molecular simulation data.

Keywords: molecular dynamics (MD); Metropolis-Hastings; Monte Carlo (MC); Markov chain Monte Carlo (MCMC); equilibration; burn-in; timeseries analysis; statistical inefficiency; integrated autocorrelation time

INTRODUCTION

Molecular simulations use Markov chain Monte Carlo (MCMC) techniques [1] to sample configurations x from an equilibrium distribution $\pi(x)$, either exactly (using Monte Carlo methods such as Metropolis-Hastings) or approximately (using molecular dynamics integrators without Metropolization) [2].

Due to the sensitivity of the equilibrium probability density $\pi(x)$ to small perturbations in configuration x and the difficulty of producing sufficiently good guesses of typical equilibrium configurations $x \sim \pi(x)$, these molecular simulations are often started from highly atypical initial conditions. For example, simulations of biopolymers might be initiated from a fully extended conformation unrepresentative of behavior in solution, or a geometry derived from a fit to diffraction data collected from a cryocooled crystal; solvated systems may be prepared by periodically replicating a small solvent box equilibrated under different conditions, yielding atypical densities and solvent structure; liquid mixtures or lipid bilayers may be constructed by using methods that fulfill spatial constraints (e.g. PackMol [3]) but create locally atypical geometries, requiring long simulation times to relax to typical configurations.

As a result, traditional practice in molecular simulation has recommended some initial portion of the trajectory be discarded to *equilibration* (also called *burn-in*¹ in the MCMC literature [4]). While the process of discarding initial samples is strictly unnecessary for the time-average of quantities of interest to eventually converge to the desired expectations [5], this nevertheless often allows the practitioner to avoid what may be impractically long run times to eliminate the bias in computed properties in finite-length simulations

induced by atypical initial starting conditions. It is worth noting that a similar procedure is not a practice universally recommended by statisticians when sampling from posterior distributions in statistical inference [4]; the differences in complexity of probability densities typically encountered in statistics and molecular simulation may explain the difference in historical practice.

As a motivating example, consider the computation of the average density of liquid argon under a given set of reduced temperature and pressure conditions shown in Figure 1. To initiate the simulation, an initial dense liquid geometry at reduced density $\rho^* \equiv \rho\sigma^3 = 0.960$ was prepared and subjected to local energy minimization. The upper panel of Figure 1 depicts the average relaxation behavior of simulations initiated from the same configuration with different random initial velocities and integrator random number seeds (see *Simulation Details*). The average of 500 realizations of this process shows a characteristic relaxation away from the initial density toward the equilibrium density (Figure 1, upper panel, black line). As a result, the expectation of the running average of the density significantly deviates from the true expectation (Figure 1, lower panel, dashed line). This effect leads to significantly biased estimates of the expectation unless simulations are sufficiently long to eliminate starting point dependent bias, which takes a surprisingly long ~ 30 ns in this example. Note that this bias is present even in the average of many realizations because the *same* atypical starting condition is used for every realization of this simulation process.

To develop an automatic approach to eliminating this bias, we take motivation from the concept of *reverse cumulative averaging* from Yang et al. [6], in which the trajectory statistics over the production region of the trajectory are examined for different choices of the end of the discarded equilibration region to determine the optimal production region to use for computing expectations and other statistical properties. We begin by first formalizing our objectives mathematically.

* Corresponding author; john.chodera@choderalab.org

¹ The term *burn-in* comes from the field of electronics, in which a short “burn-in” period is used to ensure that a device is free of faulty components—which often fail quickly—and is operating normally [4].

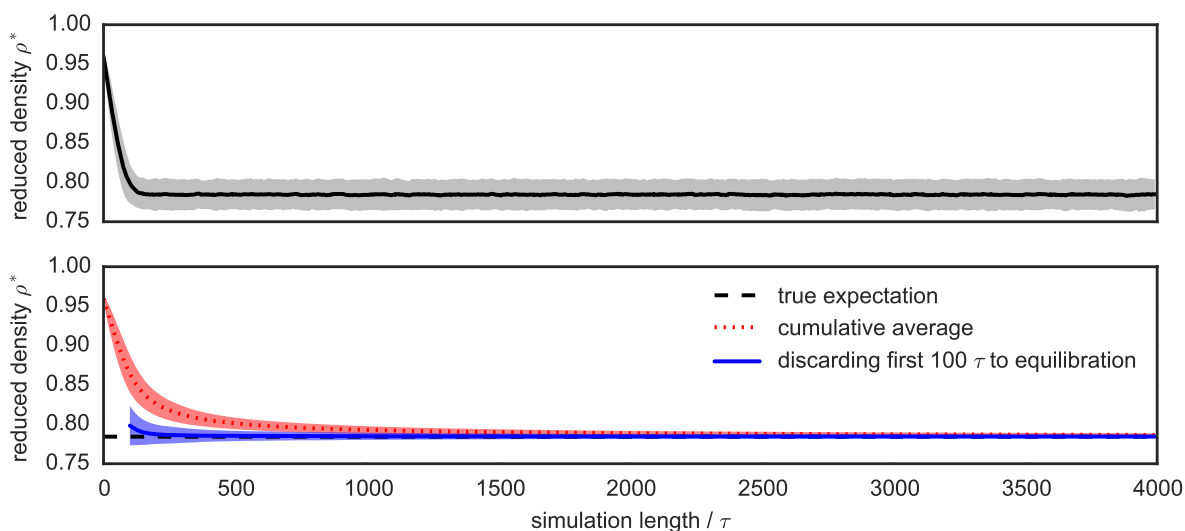


FIG. 1. Illustration of the motivation for discarding data to equilibration. To illustrate the bias in expectations induced by relaxation away from initial conditions, 500 replicates of a simulation of liquid argon were initiated from the same energy-minimized initial configuration constructed with initial reduced density $\rho^* \equiv \rho\sigma^3 = 0.960$ but different random number seeds for stochastic integration. **Top:** The average of the reduced density (black line) over the replicates relaxes to the region of typical equilibrium densities over the first $\sim 90 \tau$ of simulation time, where τ is a natural time unit (see *Simulation Details*). **Bottom:** If the average density is estimated by a cumulative average from the beginning of the simulation (red dotted line), the estimate will be heavily biased by the atypical starting density even beyond 1000τ . Discarding even a small amount of initial data—in this case 500 initial samples—results in a cumulative average estimate that converges to the true average (black dashed line) much more rapidly. Shaded regions denote 95% confidence intervals.

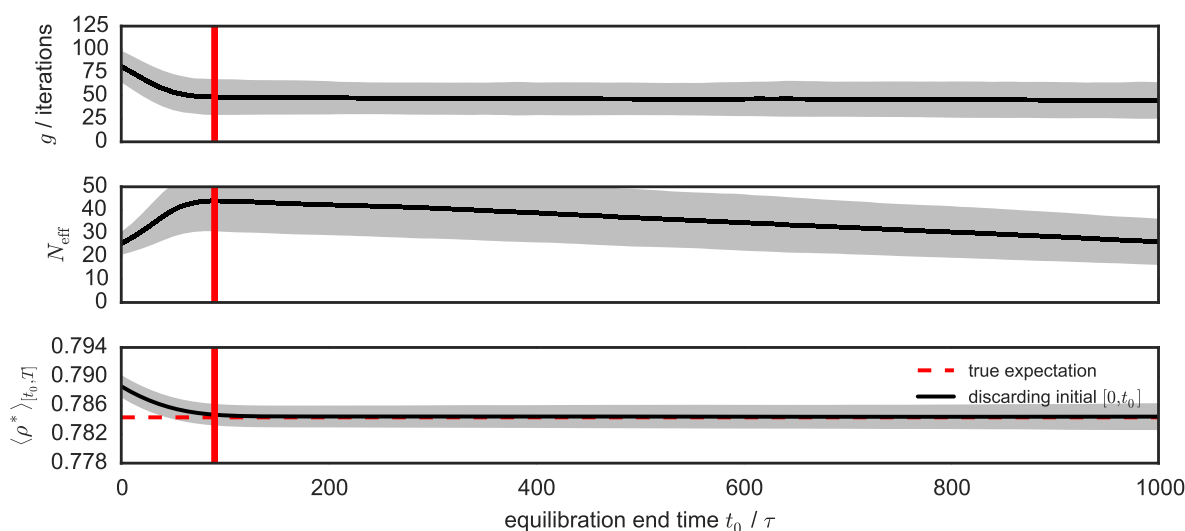


FIG. 2. Statistical inefficiency, number of uncorrelated samples, and bias for different equilibration times. Trajectories of length $T = 2000 \tau$ for the argon system described in Figure 1 were analyzed as a function of equilibration time choice t_0 . Averages over all 500 replicate simulations (all starting from the same initial conditions) are shown as dark lines, with shaded lines showing standard deviation of estimates among replicates. **Top:** The statistical inefficiency g as a function of equilibration time choice t_0 is initially very large, but diminishes rapidly after the system has relaxed to equilibrium. **Middle:** The number of effectively uncorrelated samples $N_{\text{eff}} = (T - t_0 + 1)/g$ shows a maximum at $t_0 \sim 90 \tau$ (red vertical lines), suggesting the system has equilibrated by this time. **Bottom:** The cumulative average density $\langle \rho^* \rangle_{[t_0, T]}$ computed over the span $[t_0, T]$ shows that the bias (deviation from the true estimate, shown as red dashed lines) is minimized for choices of $t_0 \geq 90 \tau$. The standard deviation among replicates (shaded region) grows with t_0 because fewer data are included in the estimate. The choice of optimal t_0 that maximizes N_{eff} (red vertical line) strikes a good balance between bias and variance. The true estimate (red dashed lines) is computed from averaging over the range $[5\,000, 10\,000] \tau$ over all 500 replicates.

88 STATEMENT OF THE PROBLEM

114 BIAS-VARIANCE TRADEOFF

89 Consider T successively sampled configurations x_t from a
90 molecular simulation, with $t = 1, \dots, T$, initiated from x_0 .
91 We presume we are interested in computing the expectation

$$\langle A \rangle \equiv \int dx A(x) \pi(x) \quad (1)$$

92 of a mechanical property of interest $A(x)$. For convenience,
93 we will refer to the timeseries $a_t \equiv A(x_t)$, with $t \in [1, T]$.
94 The estimator $\hat{A} \approx \langle A \rangle$ constructed from the entire dataset
95 is given by

$$\hat{A}_{[1, T]} \equiv \frac{1}{T} \sum_{t=1}^T a_t. \quad (2)$$

96 While $\lim_{T \rightarrow \infty} \hat{A}_{[1, T]} = \langle A \rangle$ for an infinitely long simula-
97 tion², the bias in $\hat{A}_{[1, T]}$ may be significant in a simulation of
98 finite length T .

99 By discarding samples $t < t_0$ to equilibration, we hope to
100 exclude the initial transient from our sample average, and
101 provide a less biased estimate of $\langle A \rangle$,

$$\hat{A}_{[t_0, T]} \equiv \frac{1}{T - t_0 + 1} \sum_{t=t_0}^T a_t. \quad (3)$$

102 We can quantify the overall error in an estimator $\hat{A}_{[t_0, T]}$
103 in a sample average that starts at x_0 and excludes samples
104 where $t < t_0$ by the expected error $\delta^2 \hat{A}_{[t_0, T]}$,

$$\delta^2 \hat{A}_{[t_0, T]} \equiv E_{x_0} \left[\left(\hat{A}_{[t_0, T]} - \langle A \rangle \right)^2 \right] \quad (4)$$

105 where $E_{x_0}[\cdot]$ denotes the expectation over independent re-
106 realizations of the specific simulation process initiated from
107 configuration x_0 , but with different velocities and random
108 number seeds.

109 We can rewrite the expected error $\delta^2 \hat{A}$ by separating it
110 into two components:

$$\delta^2 \hat{A}_{[t_0, T]} = E_{x_0} \left[\left(\hat{A}_{[t_0, T]} - E_{x_0}[\hat{A}_{[t_0, T]}] \right)^2 \right] + \left(E_{x_0}[\hat{A}_{[t_0, T]}] - \langle A \rangle \right)^2 \quad (5)$$

111 The first term denotes the variance in the estimator \hat{A} ,

$$\text{var}_{x_0}(\hat{A}_{[t_0, T]}) \equiv E_{x_0} \left[\hat{A}_{[t_0, T]} - E_{x_0}[\hat{A}_{[t_0, T]}] \right]^2 \quad (6)$$

112 while the second term denotes the contribution from the
113 squared bias,

$$\text{bias}_{x_0}^2(\hat{A}_{[t_0, T]}) \equiv \left(E_{x_0}[\hat{A}_{[t_0, T]}] - \langle A \rangle \right)^2 \quad (7)$$

² We note that this equality only holds for simulation schemes that sample from the true equilibrium density $\pi(x)$, such as Metropolis-Hastings Monte Carlo or Metropolized dynamical integration schemes such as hybrid Monte Carlo (HMC). Molecular dynamics simulations utilizing finite timestep integration without Metropolization will produce averages that may deviate from the true expectation $\langle A \rangle$ [2].

115 With increasing equilibration time t_0 , bias is reduced, but
116 the variance—the contribution to error due to random varia-
117 tion from having a finite number of uncorrelated samples—
118 will increase because less data is included in the estimate.
119 This can be seen in the bottom panel of Figure 2, where
120 the shaded region (95% confidence interval of the mean) in-
121 creases in width with increasing equilibration time t_0 .

122 To examine the tradeoff between bias and variance explic-
123 itly, Figure 3 plots the bias and variance (here, shown as
124 standard error) contributions against each other as a func-
125 tion of t_0 (denoted by color) as computed from statistics over
126 all 500 replicates. At $t_0 = 0$, the bias is large but variance is
127 minimized. With increasing t_0 , bias is eventually eliminated
128 but then variance rapidly grows as fewer uncorrelated sam-
129 ples are included in the estimate. There is a clear optimal
130 choice at $t_0 \sim 90 \tau$ that minimizes variance while also effec-
131 tively eliminating bias (where τ is a natural time unit—see
132 *Simulation Details*).

133 SELECTING THE EQUILIBRATION TIME

134 Is there a simple approach to choosing an optimal equi-
135 libration time t_0 that provides a significantly improved esti-
136 mate $\hat{A}_{[t_0, T]}$, even when we do not have access to multiple
137 realizations? At worst, we hope that such a procedure would
138 at least give some improvement over the naive estimate,
139 such that $\delta^2 \hat{A}_{[t_0, T]} < \delta^2 \hat{A}_{[0, T]}$; at best, we hope that we can
140 achieve a reasonable bias-variance tradeoff close to the op-
141 timal point identified in Figure 3 that minimizes bias with-
142 out greatly increasing variance. We remark that, for cases
143 in which the simulation is not long enough to reach equilib-
144 rium, no choice of t_0 will eliminate bias completely; the best
145 we can hope for is to minimize this bias.

146 While automated methods for selecting the equilibration
147 time t_0 have been proposed, these approaches have short-
148 comings that have greatly limited their use. The reverse
149 cumulative averaging (RCA) method proposed by Yang et
150 al. [6], for example, uses a statistical test for normality to de-
151 termine the point before which which the observable time-
152 series deviates from normality when examining the time-
153 series in reverse. While this concept may be reasonable
154 for experimental data, where measurements often repre-
155 sent the sum of many random variables such that the central
156 limit theorem's guarantee of asymptotic normality ensures
157 the distribution of the observable will be approximately nor-
158 mal, there is no such guarantee that instantaneous mea-
159 surements of a simulation property of interest will be nor-
160 mally distributed. In fact, many properties will be decidedly
161 *non-normal*. For a biomolecule such as a protein, for exam-
162 ple, the radius of gyration, end-to-end distance, and torsion
163 angles sampled during a simulation will all be highly non-
164 normal. Instead, we require a method that makes no as-
165 sumptions about the nature of the distribution of the prop-
166 erty under study.

AUTOCORRELATION ANALYSIS

167

168 The set of successively sampled configurations $\{x_t\}$ and
 169 their corresponding observables $\{a_t\}$ compose a correlated
 170 timeseries of observations. To estimate the statistical error
 171 or uncertainty in a stationary timeseries free of bias, we
 172 must be able to quantify the *effective number of uncorrelated*
 173 *samples* present in the dataset. This is usually accom-
 174 plished through computation of the *statistical inefficiency* g ,
 175 which quantifies the number of correlated timeseries sam-
 176 ples needed to produce a single effectively uncorrelated
 177 sample of the observable of interest. While these concepts
 178 are well-established for the analysis of both Monte Carlo
 179 and molecular dynamics simulations [7–10], we review them
 180 here for the sake of clarity.

181 For a given equilibration time choice t_0 , the statistical un-
 182 certainty in our estimator $\hat{A}_{[t_0, T]}$ can be written as,

$$\begin{aligned} \delta^2 \hat{A}_{[t_0, T]} &\equiv E_{x_0} \left[\left(\hat{A}_{[t_0, T]} - \langle \hat{A} \rangle \right)^2 \right] \\ &= E_{x_0} \left[\hat{A}_{[t_0, T]}^2 \right] - E_{x_0} \left[\hat{A}_{[t_0, T]} \right]^2 \\ &= \frac{1}{T_{t_0}^2} \sum_{t, t'=t_0}^T \{ E_{x_0} [a_t a_{t'}] - E_{x_0} [a_t] E_{x_0} [a_{t'}] \} \\ &= \frac{1}{T_{t_0}^2} \sum_{t=t_0}^T \{ E_{x_0} [a_t^2] - E_{x_0} [a_t]^2 \} \\ &\quad + \frac{1}{T_{t_0}^2} \sum_{t \neq t'=t_0}^T \{ E_{x_0} [a_t a_{t'}] - E_{x_0} [a_t] E_{x_0} [a_{t'}] \}, \end{aligned} \quad (8)$$

183 where $T_{t_0} \equiv T - t_0 + 1$, the number of correlated samples
 184 in the timeseries $\{a_t\}_{t_0}^T$. In the last step, we have split the
 185 double-sum into two separate sums—a term capturing the
 186 variance in the observations a_t , and a remaining term cap-
 187 turing the correlation between observations.

188 If t_0 is sufficiently large for the initial bias to be eliminated,
 189 the remaining timeseries $\{a_t\}_{t_0}^T$ will obey the properties of
 190 both *stationarity* and *time-reversibility*, allowing us to write,

$$\begin{aligned} \delta^2 \hat{A}_{[t_0, T]}^{\text{equil}} &= \frac{1}{T_{t_0}} [\langle a_t^2 \rangle - \langle a_t \rangle^2] \\ &\quad + \frac{2}{T_{t_0}} \sum_{n=1}^{T-t_0} \left(\frac{T_{t_0} - n}{T_{t_0}} \right) [\langle a_t a_{t+n} \rangle - \langle a_t \rangle \langle a_{t+n} \rangle] \\ &\equiv \frac{\sigma_{t_0}^2}{T_{t_0}} (1 + 2\tau_{t_0}) = \frac{\sigma_{t_0}^2}{T_{t_0}/g_{t_0}}, \end{aligned} \quad (9)$$

191 where the variance σ^2 , statistical inefficiency g , and inte-
 192 grated autocorrelation time τ (in units of the sampling in-
 193 terval) are given by

$$\sigma^2 \equiv \langle a_t^2 \rangle - \langle a_t \rangle^2, \quad (10)$$

$$\tau \equiv \sum_{t=1}^{T-1} \left(1 - \frac{t}{T} \right) C_t, \quad (11)$$

$$g \equiv 1 + 2\tau, \quad (12)$$

194 with the discrete-time normalized fluctuation autocorrela-
 195 tion function C_t defined as

$$C_t \equiv \frac{\langle a_n a_{n+t} \rangle - \langle a_n \rangle^2}{\langle a_n^2 \rangle - \langle a_n \rangle^2}. \quad (13)$$

196 In practice, it is difficult to estimate C_t for $t \sim T$, due to
 197 growth in the statistical error, so common estimators of g
 198 make use of several additional properties of C_t to provide
 199 useful estimates (see *Practical Computation of Statistical In-*
 200 *efficiencies*).

201 The t_0 subscript for the variance σ^2 , the integrated auto-
 202 correlation time τ , and the statistical inefficiency t_0 mean
 203 that these quantities are only estimated over the production
 204 portion of the timeseries, $\{a_t\}_{t=t_0}^T$. Since we assumed that
 205 the bias was eliminated by judicious choice of the equilibra-
 206 tion time t_0 , this estimate of the statistical error will be poor
 207 for choices of t_0 that are too small.

208

THE ESSENTIAL IDEA

209 Suppose we choose some arbitrary time t_0 and discard all
 210 samples $t \in [0, t_0)$ to equilibration, keeping $[t_0, T]$ as the
 211 dataset to analyze. How much data remains? We can deter-
 212 mine this by computing the statistical inefficiency g_{t_0} for
 213 the interval $[t_0, T]$, and computing the effective number of
 214 uncorrelated samples $N_{\text{eff}}(t_0) \equiv (T - t_0 + 1)/g_{t_0}$. If we
 215 start at $t_0 \equiv T$ and move t_0 to earlier and earlier points in
 216 time, we expect that the effective number of uncorrelated
 217 samples $N_{\text{eff}}(t_0)$ will continue to grow until we start to in-
 218 clude the highly atypical initial data. At that point, the inte-
 219 grated autocorrelation time τ (and hence the statistical in-
 220 efficiency g) will greatly increase (a phenomenon observed
 221 earlier, e.g. Figure 2 of [6]). As a result, the effective number
 222 of samples N_{eff} will start to plummet.

223 Figure 2 demonstrates this behavior for the liquid argon
 224 system described above, using averages of the statistical
 225 inefficiency g_{t_0} and $N_{\text{eff}}(t_0)$ computed over 500 indepen-
 226 dent replicate trajectories. At short t_0 , the average statis-
 227 tical inefficiency g (Figure 2, top panel) is large due to the
 228 contribution from slow relaxation from atypical initial con-
 229 ditions, while at long t_0 the statistical inefficiency estimate
 230 is much shorter and nearly constant of a large span of time
 231 origins. As a result, the average effective number of uncor-
 232 related samples N_{eff} (Figure 2, middle panel) has a peak at
 233 $t_0 \sim 90 \tau$ (Figure 2, vertical red lines). The effect on bias in
 234 the estimated average reduced density $\langle \rho^* \rangle$ (Figure 2, bot-
 235 tom panel) is striking—the bias is essentially eliminated for
 236 the choice of equilibration time t_0 that maximizes the num-
 237 ber of uncorrelated samples N_{eff} .

238 This suggests an alluringly simple algorithm for identify-
 239 ing the optimal equilibration time—pick the t_0 which maxi-
 240 mizes the number of uncorrelated samples N_{eff} in the time-
 241 series $\{a_t\}_{t_0}^T$ for the quantity of interest $A(x)$:

$$\begin{aligned} t_0^{\text{opt}} &= \underset{t_0}{\text{argmax}} N_{\text{eff}}(t_0) \\ &= \underset{t_0}{\text{argmax}} \frac{T - t_0 + 1}{g_{t_0}} \end{aligned} \quad (14)$$

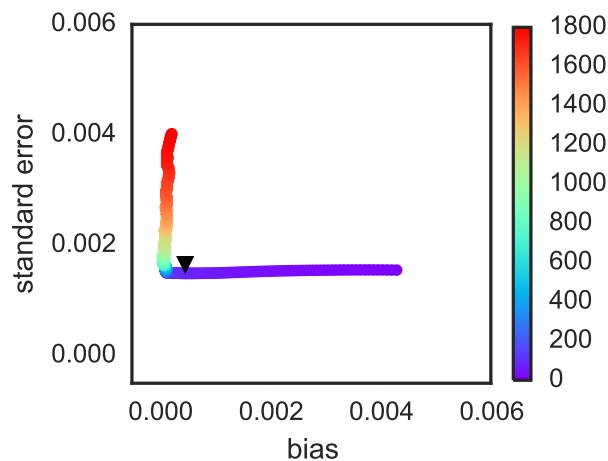


FIG. 3. Bias-variance tradeoff for fixed equilibration time versus automatic equilibration time selection. Trajectories of length $T = 2000\tau$ for the argon system described in Figure 1 were analyzed as a function of equilibration time choice t_0 , with colors denoting the value of t_0 (in units of τ) corresponding to each plotted point. Using 500 replicate simulations, the average bias (average deviation from true expectation) and standard deviation (random variation from replicate to replicate) were computed as a function of a prespecified fixed equilibration time t_0 , with colors running from violet (0τ) to red (1800τ). As is readily discerned, the bias for small t_0 is initially large, but minimized for larger t_0 . By contrast, the standard error (a measure of variance, estimated here by standard deviation among replicates) grows as t_0 grows above a certain critical time (here, $\sim 90\tau$). If the t_0 that maximizes N_{eff} is instead chosen *individually* for each trajectory based on that trajectory's estimates of statistical inefficiency $g_{[t_0, T]}$, the resulting bias-variance tradeoff (black triangle) does an excellent job minimizing bias and variance simultaneously, comparable to what is possible for a choice of equilibration time t_0 based on knowledge of the true bias and variance among many replicate estimates.

Bias-variance tradeoff. How will the simple strategy of selecting the equilibration time t_0 using Eq 14 work for cases where we do not know the statistical inefficiency g as a function of the equilibration time t_0 precisely? When all that is available is a single simulation, our best estimate of g_{t_0} is derived from that simulation alone over the span $[t_0, T]$ —will this affect the quality of our estimate of equilibration time? Empirically, this does not appear to be the case—the black triangle in Figure 3 shows the bias and variance contributions to the error for estimates computed over the 500 replicates where t_0 is individually determined from each simulation using this simple scheme based on selecting t_0 to maximize N_{eff} for each individual realization. Despite not having knowledge about multiple realizations, this strategy effectively achieves a near-optimal balance between minimizing bias without increasing variance.

Overall RMS error. How well does this strategy perform in terms of decreasing the *overall* error $\delta\hat{A}_{[t_0, T]}$ compared to $\delta\hat{A}_{[0, T]}$? Figure 4 compares the expected standard error (denoted $\delta\hat{A}$) as a function of a fixed initial equilibration time t_0 (black line with shaded region denoting 95% confidence interval) with the strategy of selecting t_0 to maximize N_{eff} for each realization (red line with shaded region denoting 95% confidence interval). While the minimum error for the fixed- t_0 strategy (0.00154 ± 0.00005) is achieved at 90τ —a fact that could only be determined from knowledge of multiple realizations—the simple strategy of selecting t_0 using Eq. 14 achieves a minimum error of 0.00171 ± 0.00006 , only 11% worse (compared to errors of 0.00456 ± 0.00007 , or 296% worse, should no data have been discarded).

DISCUSSION

The scheme described here—in which the equilibration time t_0 is computed using Eq. 14 as the choice that maximizes the number of uncorrelated samples in the production region $[t_0, T]$ —is both conceptually and computationally straightforward. It provides an approach to determining the optimal amount of initial data to discard to equilibrate in order to minimize variance while also minimizing initial bias, and does this without employing statistical tests that require generally unsatisfiable assumptions of normality of the observable of interest. All that is needed is to save the timeseries $\{a_t\}_1^T$ of the observable $A(x)$ of interest—there is no need to store full configurations x_t —and post-process this dataset with a simple analysis procedure, for which we have provided a convenient Python reference implementation (see *Simulation Details*). As we have seen, this scheme empirically appears to select a practical compromise between bias and variance even when the statistical inefficiency g is estimated directly from the trajectory using Eq. 12.

To show that this approach is indeed general, we repeated the analysis illustrated above in Figs. 1–4 for a different choice of observable $A(x)$ for the same liquid argon system—in this case, the reduced potential energy³ $u^*(x) \equiv \beta U(x)$. The results of this analysis are collected in Fig. 5. As can readily be seen, this reduced potential behaves in essentially the same way the reduced density does, and the simple scheme for automated determination of equilibration time t_0 from Eq. 14 does just as well.

A word of caution is necessary. One can certainly envision pathological scenarios where this algorithm for selecting an optimal equilibration time will break down. In cases where the simulation is not long enough to reach equilibrium—let

³ Note that the *reduced potential*[11] for the isothermal-isobaric ensemble is generally defined as $u^*(x) = -\ln \pi(x) = \beta[u(x) + pV(x)]$ to include the pressure-volume term $\beta pV(x)$, but in order to demonstrate the performance of this analysis on an observable distinct from the density, we omit the pressure-volume term in the present analysis.

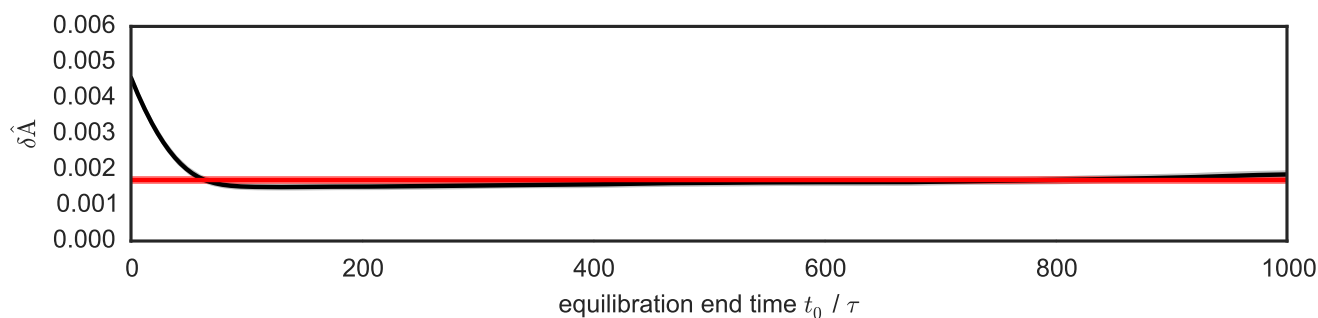


FIG. 4. RMS error for fixed equilibration time versus automatic equilibration time selection. Trajectories of length $T = 2000\tau$ for the argon system described in Figure 1 were analyzed as a function of fixed equilibration time choice t_0 . Using 500 replicate simulations, the root-mean-squared (RMS) error (Eq. 4) was computed (black line) along with 95% confidence interval (gray shading). The RMS error is minimized for fixed equilibration time choices in the range 90–200 τ . If the t_0 that maximizes N_{eff} is instead chosen *individually* for each trajectory based on that trajectory’s estimated statistical inefficiency $g_{[t_0, T]}$ using Eq. 14, the resulting RMS error (red line, 95% confidence interval shown as red shading) is quite close to the minimum RMS error achieved from any particular *fixed* choice of equilibration time t_0 , suggesting that this simple automated approach to selecting t_0 achieves close to optimal performance.

305 alone collect many uncorrelated samples from it—no choice
 306 of equilibration time will bestow upon the experimenter the
 307 ability to produce an unbiased estimate of the true expecta-
 308 tion. Similarly, in cases where insufficient data is available
 309 for the statistical inefficiency to be estimated well, this al-
 310 gorithm is expected to perform poorly. However, in these
 311 cases, the data itself should be suspect if the trajectory is not
 312 at least an order of magnitude longer than the minimum es-
 313 timated autocorrelation time.

314 SIMULATION DETAILS

315 All molecular dynamics simulations described here were
 316 performed with OpenMM 6.2 [12] (available at openmm.org)
 317 using the Python API. All scripts used to retrieve the software
 318 versions used here, run the simulations, analyze data, and
 319 generate plots—along with the simulation data itself and
 320 scripts for generating figures—are available on GitHub⁴.

321 To model liquid argon, the LennardJonesFluid model
 322 system in the [openmmtools](http://openmmtools.org) package⁵ was used with param-
 323 eters appropriate for liquid argon ($\sigma = 3.4 \text{ \AA}$, $\epsilon = 0.238$
 324 kcal/mol). All results are reported in reduced (dimension-
 325 less) units. A cubic switching function was employed, with
 326 the potential gently switched to zero over $r \in [\sigma, 3\sigma]$, and
 327 a long-range isotropic dispersion correction accounting for
 328 this switching behavior used to include neglected contribu-
 329 tions. Simulations were performed using a periodic box of
 330 $N = 500$ atoms at reduced temperature $T^* \equiv k_B T / \epsilon =$
 331 0.850 and reduced pressure $p^* \equiv p\sigma^3 / \epsilon = 1.266$ using a

⁴ All Python scripts necessary to reproduce this work—along with data plotted in the published version—are available at:

<http://github.com/choderalab/automatic-equilibration-detection>

⁵ available at <http://github.com/choderalab/openmmtools>

332 Langevin integrator [13] with timestep $\Delta t = 0.01\tau$ and col-
 333 lision rate $\nu = \tau^{-1}$, with characteristic oscillation timescale
 334 $\tau = \sqrt{mr_0^2 / 72\epsilon}$ and $r_0 = 2^{1/6}\sigma$ [14]. All times are reported
 335 in multiples of the characteristic timescale τ . A molecu-
 336 lar scaling Metropolis Monte Carlo barostat with Gaussian
 337 simulation volume change proposal moves attempted every
 338 τ (100 timesteps), using an adaptive algorithm that ad-
 339 justs the proposal width during the initial part of the simu-
 340 lation [12]. Densities were recorded every τ (100 timesteps).
 341 The true expectation $\langle \rho^* \rangle$ was estimated from the sample
 342 average over all 500 realizations over [5000, 10000] τ .

343 The automated equilibration detection scheme is also
 344 available in the `timeseries` module of the `pymbar` pack-
 345 age as `detectEquilibration()`, and can be accessed us-
 346 ing the following code:

```

347 from pymbar.timeseries import detectEquilibration
348 # determine equilibrated region
349 [t0, g, Neff_max] = detectEquilibration(A_t)
350 # discard initial samples to equilibration
351 A_t = A_t[t0:]
    
```

347 PRACTICAL COMPUTATION OF STATISTICAL INEFFICIENCIES

348 The robust computation of the statistical inefficiency g
 349 (defined by Eq. 12) for a finite timeseries $a_t, t = 0, \dots, T$
 350 deserves some comment. There are, in fact, a variety of
 351 schemes for estimating g described in the literature, and
 352 their behaviors for finite datasets may differ, leading to dif-
 353 ferent estimates of the equilibration time t_0 using the algo-
 354 rithm of Eq. 14.

355 The main issue is that a straightforward approach to es-
 356 timating the statistical inefficiency using Eqs. 11–13 in which
 357 the expectations are simply replaced with sample estimates
 358 causes the statistical error in the estimated correlation func-

tion C_t to grow with t in a manner that allows this error to quickly overwhelm the sum of Eq. 11. As a result, a number of alternative schemes—generally based on controlling the error in the estimated C_t or truncating the sum of Eq. 11 when the error grows too large—have been proposed.

For stationary, irreducible, reversible Markov chains, Geyer observed that a function $\Gamma_k \equiv \gamma_{2k} + \gamma_{2k+1}$ of the unnormalized fluctuation autocorrelation function $\gamma_t \equiv \langle a_i a_{i+t} \rangle - \langle a_i \rangle^2$ has a number of pleasant properties (Theorem 3.1 of [15]): It is strictly positive, strictly decreasing, and strictly convex. Some or all of these properties can be exploited to define a family of estimators called *initial sequence methods* (see Section 3.3 of [15] and Section 1.10.2 of [4]), of which the *initial convex sequence* (ICS) estimator is generally agreed to be optimal, if somewhat more complex to implement.⁶

All computations in this manuscript used the fast multiscale method described in Section 5.2 of [10], which we found performed equivalently well to the Geyer estimators (data not shown). This method is related to a multiscale variant of the *initial positive sequence* (IPS) method of Geyer [16], where contributions are accumulated at increasingly longer lag times and the sum of Eq. 11 is truncated when the terms become negative. We have found this method to be both fast

and to provide useful estimates of the statistical inefficiency, but it may not perform well for all problems.

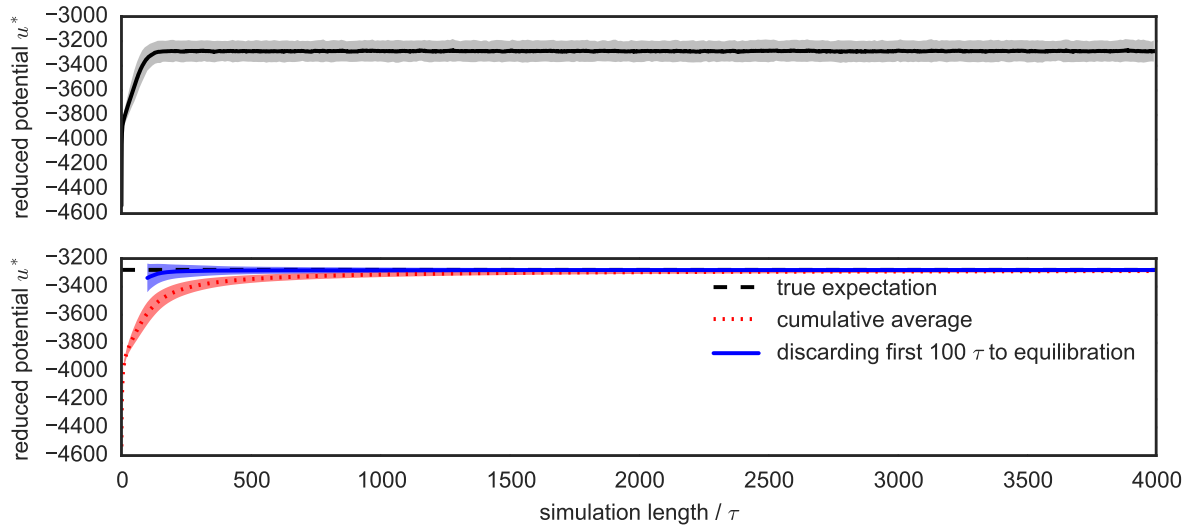
ACKNOWLEDGMENTS

We are grateful to William C. Swope (IBM Almaden Research Center) for his illuminating introduction to the use of autocorrelation analysis for the characterization of statistical error, as well as Michael R. Shirts (University of Virginia), David L. Mobley (University of California, Irvine), Michael K. Gilson (University of California, San Diego), Kyle A. Beauchamp (MSKCC), and Robert C. McGibbon (Stanford University) for valuable discussions on this topic, and Joshua L. Adelman (University of Pittsburgh) for helpful feedback and encouragement. We are grateful to Michael K. Gilson (University of California, San Diego), Wei Yang (Florida State University), Sabine Reißer (SISSA, Italy), and the anonymous referees for critical feedback on the manuscript itself. JDC acknowledges a Louis V. Gerstner Young Investigator Award, NIH core grant P30-CA008748, and the Sloan Kettering Institute for funding during the course of this work.

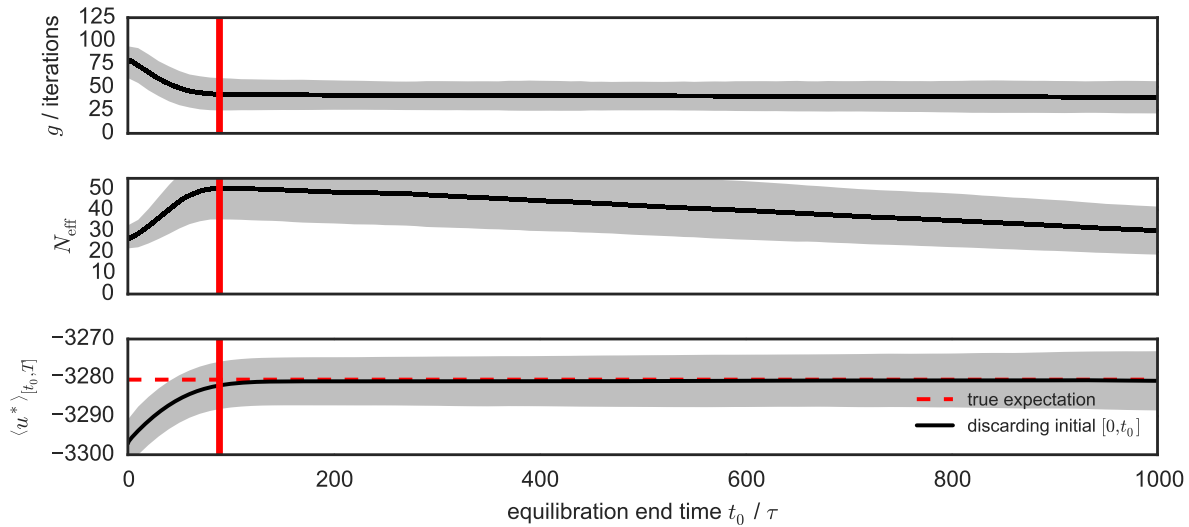
-
- [1] J. S. Liu, *Monte Carlo strategies in scientific computing*, 2nd ed. (Springer-Verlag, New York, 2002).
- [2] D. Sivak, J. Chodera, and G. Crooks, *Physical Review X* **3**, 011007 (2013), bibtext: Sivak:2013:Phys.Rev.X.
- [3] L. Martínez, R. Andrade, E. G. Birgin, and J. M. Martínez, *J. Chem. Theor. Comput.* **30**, 2157 (2009).
- [4] S. Brooks, A. Gelman, G. L. Jones, and X.-L. Meng, in *Handbook of Markov chain Monte Carlo*, Chapman & Hall/CRC *Handbooks of Modern Statistical Methods* (CRC Press, ADDRESS, 2011), Chap. Introduction to Markov chain Monte Carlo.
- [5] C. Geyer, Burn-in is unnecessary., <http://users.stat.umn.edu/~geyer/mcmc/burn.html>.
- [6] W. Yang, R. Bittetti-Putzer, and M. Karplus, *J. Chem. Phys.* **120**, 2618 (2004).
- [7] H. Müller-Krumbhaar and K. Binder, *J. Stat. Phys.* **8**, 1 (1973).
- [8] W. C. Swope, H. C. Andersen, P. H. Berens, and K. R. Wilson, *J. Chem. Phys.* **76**, 637 (1982).
- [9] W. Janke, in *Quantum Simulations of Complex Many-Body Systems: From Theory to Algorithms*, edited by J. Grotendorst, D. Marx, and A. Murmatsu (John von Neumann Institute for Computing, ADDRESS, 2002), Vol. 10, pp. 423–445.
- [10] J. D. Chodera, W. C. Swope, J. W. Pitera, C. Seok, and K. A. Dill, *J. Chem. Theor. Comput.* **3**, 26 (2007).
- [11] M. R. Shirts and J. D. Chodera, *J. Chem. Phys.* **In press**, (2008).
- [12] P. Eastman, M. Friedrichs, J. D. Chodera, R. Radmer, C. Bruns, J. Ku, K. Beauchamp, T. J. Lane, L.-P. Wang, D. Shukla, T. Tye, M. Houston, T. Stitch, and C. Klein, *J. Chem. Theor. Comput.* **9**, 461 (2012).
- [13] D. A. Sivak, J. D. Chodera, and G. E. Crooks, *J. Phys. Chem. B* **118**, 6466 (2014).
- [14] B. Veytsman and M. Kotelyanskii, Lennard-Jones potential revisited., <http://borisv.lk.net/matsc597c-1997/simulations/Lecture5/node3.html>.
- [15] C. J. Geyer, *Stat. Sci.* **76**, 473 (1992).
- [16] C. J. Geyer and E. A. Thompson, *J. Royal Stat. Soc. B* **54**, 657 (1992).

⁶ Implementations of these methods are provided with the code dis-

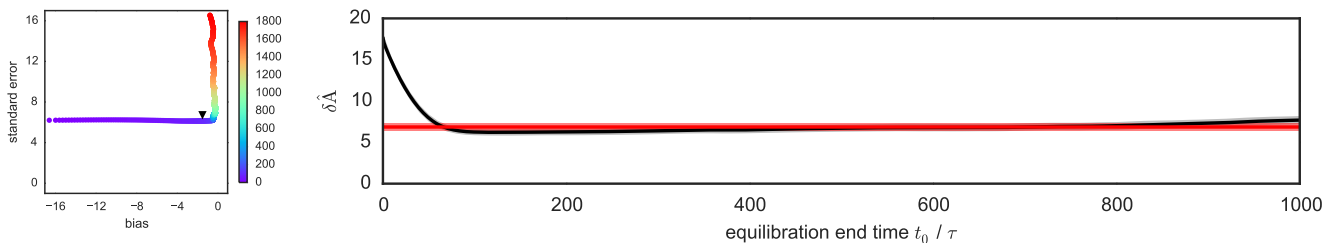
tributed with this manuscript.



(a) *Top*: Average reduced potential for relaxation from initial conditions (black solid line). *Bottom*: Cumulative averages of reduced potential with (blue solid line) and without (red dotted line) discarding some initial data to equilibration.



(b) Statistical inefficiency g , effective number of uncorrelated samples N_{eff} , and average of reduced potential u^* omitting equilibration time t_0 for analysis of reduced potential.



(c) Bias-variance tradeoff for reduced potential. (d) Overall RMSE error in estimate of reduced potential as a function of discarded equilibration time (black) compared with RMSE from automatic equilibration detection scheme (red line).

FIG. 5. Corresponding analysis for reduced potential energy of liquid argon system. The analyses of Figs. 1–4 were repeated for the reduced potential energy $u^*(x) \equiv \beta U(x)$ of the liquid argon system. As with the analysis of reduced density, the simple automated determination of equilibration time t_0 from Eq. 14 works equivalently well for the reduced potential. Shaded regions denote 95% confidence interval.