

# Building Markov State Models for Periodically Driven Non-Equilibrium Systems

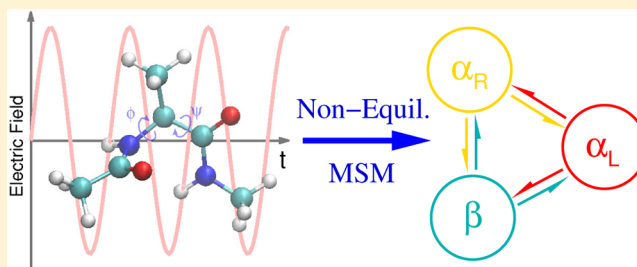
Han Wang<sup>\*,†,‡</sup> and Christof Schütte<sup>\*,‡,¶</sup>

<sup>†</sup>CAEP Software Center for High Performance Numerical Simulation, Beijing, China

<sup>‡</sup>Zuse Institute Berlin (ZIB), Berlin, Germany

<sup>¶</sup>Institute for Mathematics, Freie Universität Berlin, Berlin, Germany

**ABSTRACT:** Non-equilibrium molecular dynamics (NEMD) simulations have seen increased interest the past few years, especially for molecular systems with periodic forcing by external fields (e.g., in the context of studying the effects of electromagnetic radiation on human body tissues). Recently, an NEMD method with local thermostating has been proposed that allows for the study of non-equilibrium processes in a statistically reliable and thermodynamically consistent way. In this article, we demonstrate how to construct Markov state models (MSMs) for such NEMD simulations. MSM building is well-established for systems in equilibrium, where MSMs with only a few (macro-)states allow for accurate reproduction of the essential kinetics of the molecular system under consideration. Non-equilibrium MSMs have not yet been established. This article presents a method for constructing such MSMs and illustrates their validity and usefulness through conformation dynamics of an alanine dipeptide in an external electric field.



## 1. INTRODUCTION

Biomolecular systems under non-equilibrium conditions caused by external fields, especially systems under periodic forcing, have attracted increasing interest recently. For example, the potential effects of electromagnetic radiation on human tissues (e.g., DNA, protein, and membranes) have been extensively investigated in a vast number of studies with the following references representing an incomplete selection.<sup>1–14</sup> Molecular dynamics (MD) simulations have proven particularly useful for understanding the response of biomolecular conformations to external fields because of their ability to resolve molecular details that sometimes cannot be resolved through experiments. Only recently has a non-equilibrium MD simulation (NEMD) method with local thermostating been proposed<sup>15</sup> that allows for studying non-equilibrium processes in a statistically reliable and thermodynamically consistent way. Despite the significance of non-equilibrium phenomena, analysis of the non-equilibrium MD simulations mainly follows standard approaches, and reliable tools for the quantitative description of the essential conformational dynamics of molecular systems under external forcing remain unavailable.

Despite their many advantages, MD simulations have severe limitations. For example, one has to assume that the underlying force fields describe internal and external molecular interactions appropriately, and the maximum possible simulation length often is shorter than the time scale of interest. This article is primarily focused on circumventing the latter obstacle by introducing non-equilibrium Markov state models (MSMs). Over the past decade, MSMs have been thoroughly developed in theory,<sup>16,17</sup> applications (see ref 18 for a recent overview), and software

implementations<sup>19,20</sup> but only for systems under equilibrium conditions. The main concept underlying equilibrium MSMs is to approximate the original high-dimensional MD system using reduced Markovian dynamics over a finite number of (macro-)states. These (macro-)states have to be identified with the dominant metastable sets because typical MD trajectories remain close to a metastable set substantially longer than necessary for transition to another such set.<sup>16,21</sup> In this case, the metastable sets are the main conformations of the molecular system under consideration, which often enough are given by the main wells in the energy landscape. It has been shown that for molecular systems exhibiting such metastable sets, the Markovian dynamics given by an MSM allow very close approximation of the longest relaxation processes of the underlying molecular system, at least under equilibrium conditions.<sup>22,23</sup> Moreover, it has been demonstrated that, in such cases, MSM building requires only short MD trajectories, which are much shorter than the time scales of interest.<sup>24,25</sup> Thus, MSM building often allows the study of dynamic behavior on long time scales without requiring MD trajectories of comparable length. Moreover, MSMs are utilized for understanding very long MD simulations. Extracting the essential structures and dynamic properties from long MD runs is becoming increasingly difficult as the system sizes and trajectory lengths grow. MSMs have been used to construct kinetic fingerprints from MD simulations,<sup>26</sup> which facilitate understanding of essential dynamics and comparisons with experimental data.<sup>27</sup>

Received: November 6, 2014

Published: February 25, 2015

To the best of our knowledge, this work is the first to attempt using MSMs to analyze non-equilibrium systems under periodic external forcing. More precisely, we will demonstrate how to use MSMs to investigate the conformational dynamics of a peptide (alanine dipeptide) under an oscillating electric field (EF). To this end, we will show how to generalize Markov state modeling to periodic non-equilibrium conditions, where reversibility of the dynamics cannot be assumed as is often the case in literature on MSM building.

The article is outlined as follows. In section 2, we discuss the temporal and spatial discretizations needed to construct an MSM. We consider spatial discretization of the dihedral angle space in the traditional sense of full partition MSMs.<sup>28,29</sup> In the temporal direction, the non-equilibrium process is discretized utilizing Floquet's theorem. This results in a Markov process that is time-homogeneous but not necessarily reversible, that is, irreversible. Because a full spatial partition is not rational for high-dimensional systems, we show in section 3 how to construct a few-states MSM based on milestoning<sup>21,29</sup> of the discretized irreversible Markov process. The discretizations and resulting MSM are validated in section 4 by comparing the kinetic fingerprints given by the MSM to brute force NEMD simulations of alanine dipeptide under an oscillating EF. The findings are summarized in section 5, which includes a chart presenting the workflow for MSM building of non-equilibrium systems as well as a list of open questions.

## 2. NON-EQUILIBRIUM MOLECULAR DYNAMICS AND ITS DISCRETIZATION

We consider diffusive molecular dynamics in an energy landscape  $V$  driven by the time-dependent external driving force  $E(t)D(x_t)$  with the  $T$ -periodic external field  $E(t)$

$$dx_t = (-\nabla V(x_t) + E(t)D(x_t)) dt + \sqrt{2\beta^{-1}} dw_t \quad (1)$$

where  $x_t \in \Omega$  denotes the state of the molecular system at time  $t$  in state space  $\Omega \subset \mathbb{R}^{3N_{\text{atom}}}$  ( $N_{\text{atom}}$  being the number of atoms),  $w_t$  denotes standard  $3N_{\text{atom}}$ -dimensional Brownian motion, and  $\beta$  denotes the inverse temperature such that  $\beta = 1/(k_B T)$ . Thermostated Hamiltonian or Langevin dynamics can be treated in the same way as explained herein; thus, for the sake of simplicity, we focus on the discussion of diffusive dynamics. The propagation of probability densities  $\rho = \rho(x, t)$ , based on this kind of dynamics in the sense of  $\rho(x, t) dx = P[x_t \in [x, x + dx]]$ , is governed by the Fokker–Planck equation

$$\frac{\partial \rho}{\partial t} = \mathcal{L}^\dagger(t)\rho \quad (2)$$

where  $\mathcal{L}^\dagger(t)$  is the adjoint of the generator

$$\mathcal{L}(t) = \beta^{-1}\Delta_x + (-\nabla_x V(x) + E(t)D(x))\nabla_x \quad (3)$$

where  $\Delta_x$  denotes the Laplacian operator and  $\nabla_x$  denotes the nabla operator with respect to  $x$ . The periodicity of the external driving force induces the periodicity of the generator such that  $\mathcal{L}(t) = \mathcal{L}(t + T)$ .

**2.1. Spatial Discretization: Master Equation.** We will now introduce an appropriate spatial discretization for this kind of NEMD, which is being done for reasons of simplicity only; it could be completely avoided at the price of more technical arguments. For achieving this discretization, we introduce a partition of state space  $\Omega$  into a finite number of disjoint sets

$\{\Omega_1, \dots, \Omega_N\}$ , satisfying  $\Omega = \cup_i \Omega_i$ ,  $\Omega_i \cap \Omega_j = \emptyset$ ,  $\forall i \neq j$ . Utilizing the procedure described in ref 30, the original Fokker–Planck in eq 2 is discretized, resulting in a time-inhomogeneous Markov jump process in state space  $S = \{1, \dots, N\}$  with a time-dependent rate matrix  $L(t) \in \mathbb{R}^{N \times N}$  satisfying

$$\sum_{j=1}^N L_{ij}(t) = 0 \quad (4)$$

$$L_{ij}(t) \geq 0, \quad i \neq j \quad (5)$$

$$L_{ij}(t) = L_{ij}(t + T) \quad (6)$$

for all real time  $t \geq 0$ . Moreover, the rate matrix  $L$  has the form  $L(t) = L_0 + E(t)L_1$ , where  $E(t)$  is periodic with period  $T > 0$ . Analogous to eq 2, the Markov jump process generated by  $L(t)$  transports probability distributions according to the associated Master equation

$$\frac{dp(t)}{dt} = L^T(t)p(t) \quad (7)$$

where  $L^T(t)$  denotes the matrix transpose of  $L(t)$  and  $p(t)$  is an  $N$ -vector denoting the probability distribution on  $S$  at time  $t$ ;  $p(i, t)$ , for example, is the probability of being in state  $i$  (which corresponds to set  $\Omega_i$ ) at time  $t$ . As usual, properties 4 and 5 of  $L(t)$  guarantee that the total probability mass is conserved (i.e., if  $p(i, 0) \geq 0$  component-wise, then  $p(i, t) \geq 0$  and  $\sum_i p(i, t) = \sum_i p(i, 0)$ ). The temporal evolution of the probability distribution  $p(t)$  can be formally written as

$$p(t) = \Phi(t)p(0) \quad (8)$$

by using the associated propagator matrix  $\Phi(t) \in \mathbb{R}^{N \times N}$  that solves

$$\frac{d}{dt}\Phi(t) = L^T(t)\Phi(t), \quad \Phi(0) = \text{Id} \quad (9)$$

Because the last equation can be considered column-wise, the propagator matrix inherits column-wise conservation properties:  $\Phi_{ij}(t) \geq 0$  and  $\sum_{i=1}^N \Phi_{ij}(t) = 1$ . That is,  $\Phi^T(t)$  is a stochastic matrix satisfying  $\Phi^T(t)e = e$ , with  $e = (1, \dots, 1)^T \in \mathbb{R}^N$ . Regarding these considerations, we find

$$\Phi_{ij}(t) = P(X_t = i | X_0 = j) \quad (10)$$

where  $X_t$  denotes the Markov process generated by  $L(t)$ .

The discretization sets that we used to go from  $x_t$  and  $\mathcal{L}(t)$  to  $X_t$  and  $L(t)$ , respectively, can be assumed to provide an arbitrarily fine partitioning of the original state space; then, the transport properties of  $L(t)$  are almost perfect approximations of the transport properties of  $\mathcal{L}(t)$ , in particular, the approximation  $p(i, t) \approx P(x_t \in \Omega_i)$  is almost perfect.

**2.2. Temporal Discretization: Floquet Theorem.** As an effect of the periodicity of  $L(t)$ , the propagator  $\Phi(t + T)$  satisfies

$$\Phi(t + T) = \Phi(t)\Phi(T) \quad (11)$$

for all  $t \geq 0$ . This can be seen by considering  $Y(t) = \Phi(t + T)$ . It satisfies

$$\frac{d}{dt}Y(t) = L^T(t + T)Y(t) = L^T(t)Y(t), \quad Y(0) = \Phi(T)$$

When we consider this identity column-wise and use the propagator property of  $\Phi(t)$ , we get  $\Phi(t + T) = Y(t) = \Phi(t)\Phi(T)$ . As a consequence of eq 11, we get for all integers  $m = 0, 1, 2, \dots$  that

$$\Phi(t + mT) = \Phi(t)\Phi^m(T) \quad (12)$$

In combination with eq 8, we therefore know the solution  $p(t)$  of the Master equation for all  $t \geq 0$  if we can compute  $\Phi(t)$  for  $t \in (0, T)$ . This is known as the Floquet theorem.<sup>31</sup> In particular, we get the long-term evolution of the propagator

$$\Phi(mT) = \Phi^m(T) \quad (13)$$

where  $\Phi^m(T)$  denotes the  $m$ th power of  $\Phi(T)$ . Thus, for the probability at integral periods, we have

$$p(mT) = \Phi(mT)p(0) = \Phi^m(T)p(0) \quad (14)$$

Using the Floquet theorem, the time-inhomogeneous Markov jump process  $X_t$  is therefore discretized into a time-homogeneous (but not necessarily reversible) Markov jump process  $\tilde{X}_m = X_{mT}$ ,  $m \in \mathbb{N}$ , which is generated by the transition matrix

$$P = \Phi^T(T) \quad (15)$$

We prefer to consider the discrete-time process  $\tilde{X}_m$  instead of the time-continuous process  $X_t$  because the powerful theories and computational tools for time-homogeneous Markov processes can be applied directly. It is worth noting that many of these tools require the transition matrix  $P$  to satisfy the detailed balance condition. The computations in Appendix A show that  $P$  will in general not satisfy this condition; in fact, deviation from reversibility can be estimated from the work periodic driving does to the system. There is no doubt that information within one period is lost by using this temporal discretization; however, information regarding the long-term behavior of the system on time scales much longer than the period will be described perfectly because of  $\tilde{X}_m = X_{mT} \approx x_{mT}$  whenever our spatial discretization is sufficiently fine. At the same time, the computational cost of generating  $\tilde{X}_m$  is much less demanding than the brute force simulations of NEMD, which implies lower statistical uncertainty in calculating the observables of interest.

Because  $P$  is a stochastic matrix, its eigenvalues are contained in the unit circle in the complex plane (i.e., each eigenvalue  $\lambda$  (potentially complex-valued) satisfies  $|\lambda| \leq 1$ ). Furthermore,  $\lambda = 1$  is an eigenvalue with a right eigenvector  $e = (1, \dots, 1)^T$  and a left eigenvector  $\mu$  satisfying  $\mu^T P = \mu^T$ . From now on, we assume  $P$  to be irreducible and aperiodic, such that the Perron–Frobenius theorem holds; thus, the eigenvector corresponding to the eigenvalue  $\lambda = 1$  is non-negative component-wise, and unique (up to normalization  $\sum_i \mu(i) = 1$ ). In this case,  $\mu$  is the stationary measure in the sense that  $\mu^T P^m = \mu^T$ ,  $m \in \mathbb{N}$ , and more precisely, the asymptotic evolution of an initial probability distribution  $p(t = 0)$  by the process satisfies  $p^T(0)P^m \rightarrow \mu^T$ ,  $m \rightarrow \infty$ , so that  $\mu$  can be seen as the quasi-stationary distribution of the nonstationary process.

### 3. MARKOV STATE MODEL

If the discretization cells  $\Omega_i$ ,  $i = 1, \dots, N$ , form a fine partition of the molecular state space, the Markov chain defined via the transition matrix  $P$  is discrete in time, but in space it remains a fine-scale description of the transport properties of the dynamics with a very large number ( $N$ ) of states. Now, we want to coarse our

description much further by constructing a Markov state model (MSM) for  $P$  with  $K \ll N$  macrostates, which would be the metastable states of the system. The resulting  $K \times K$  MSM transition matrix  $\hat{P}$  then defines the coarse-grained long-term kinetics that shall approximate the original long-term kinetics well. The idea behind MSM building is that, given that the molecular system under consideration exhibits metastable conformations, it is usually possible to construct a relatively small number of discrete sets—the metastable sets that form the so-called macrostates—that correctly describe the slow dynamics, and in each set, the fast dynamics relaxes on time scales that are significantly shorter than the metastable time scales. Then, if the MSM dynamics reproduces the slow time scales and the corresponding transitions of the original dynamics, the former is considered to be a good approximation of the latter.

MSM building has attracted much attention recently, and extensive theories,<sup>16</sup> algorithms,<sup>18</sup> applications (for examples, see refs 24 and 32 for two of hundreds of articles), and software<sup>19,20</sup> have been developed. However, by far most of the literature is related to building standard MSMs for equilibrium MD. In standard MSMs, the transition region also has to be discretized, a feature that often forces the user to incorporate more macrostates than is essentially needed to approximate long-term kinetics. In refs 16, 21, 22, and 33, it is shown how to construct a nonstandard MSM in a way that avoids this problem for equilibrium MD, that is, if  $P$  satisfies the detailed balance condition of (1) identifying the cores of the metastable sets of the dynamics and (2) uses them as milestones to construct an MSM in which the macrostates are the metastable core sets and  $\hat{P}$  is the transition matrix of the milestone process<sup>16,21,29</sup> that models the jumping behavior of the original dynamics between the metastable regions.

However, because we cannot assume  $P$  to satisfy such a detailed balance, we instead follow the approach for nonstandard MSM construction recently proposed in ref 34, which allows identification of the metastable core sets for the nonreversible transition matrix  $P$ . Assume that this approach leads to the  $K$  core sets  $C_1, \dots, C_K \subset S$  that are appropriate metastable sets. Following refs 16 and 21, the process  $(\tilde{X}_m)$  associated with  $P$  is coarse-grained into the so-called milestone process  $(\hat{X}_m)$  in the following way:  $\hat{X}_m$  only has  $K$  states associated with the sets  $C_j$ ,  $j = 1, \dots, K$ . The sequence of random variables  $(\hat{X}_m)$  is defined via the sequence  $(\tilde{X}_m)$  (i.e., trajectories of  $(\tilde{X}_m)$  induce trajectories of  $(\hat{X}_m)$ ). We set  $\hat{X}_m = j$  if the last core set that the process  $(\tilde{X}_m)$  entered prior to or at time  $m$  is the core set  $C_j$ .

Now, consider an arbitrary infinitely long trajectory of  $(\tilde{X}_m)$ . Because of ergodicity, we know that the states in this trajectory will be distributed according to the quasi-stationary distribution  $\mu$ . On the basis of such an infinitely long trajectory, we can consider the probability  $q_j^-(i)$  that conditioned on  $\tilde{X}_m = i$  the last core set hit has been  $C_j$ . This function is called the backward committor of  $(\tilde{X}_m)$  associated with the set  $C_j$  and is associated with the milestone process via

$$q_j^-(i) = P_\mu(\hat{X}_m = j | \tilde{X}_m = i) \quad (16)$$

where the index  $\mu$  refers to the fact that  $\hat{X}_m$  is distributed due to  $\mu$ . From the last equation, the stationary distribution of the milestone process is given by

$$\hat{\mu}_j = \sum_{i \in S} q_j^-(i) \mu(i) \quad (17)$$



that is, the probability to find  $\tilde{X}_m = j$  in infinitely long trajectories of the milestoneoning process is  $\hat{\mu}_j$ .

Next, one also has to consider the forward committor  $q_j^+(i)$  identical to the probability that conditioned on  $\tilde{X}_m = i$ , such that the next core set to be hit will be  $C_j$ . The forward and backward committors  $q_j^+$  and  $q_j^-$  for each core set  $C_j$  can be computed from  $P$  by solving the linear equations<sup>35</sup>

$$\begin{aligned}(P - \text{Id})q_j^+(i) &= 0, \quad i \in C \\ q_j^+(i) &= 1, \quad i \in C_j \\ q_j^+(i) &= 0, \quad i \in C_k, k \neq j\end{aligned}\quad (18)$$

where  $C = S \setminus \cup_j C_j$  and

$$\begin{aligned}(P^b - \text{Id})q_j^-(i) &= 0, \quad i \in C \\ q_j^-(i) &= 1, \quad i \in C_j \\ q_j^-(i) &= 0, \quad i \in C_k, k \neq j\end{aligned}\quad (19)$$

where  $P^b$  denotes the transition matrix of the time-reversed process given by  $P_{ji}^b = \mu(i)P_{ij}/\mu(j)$ .

We define the one-step transition matrix  $\hat{P}$  for the milestone process by

$$\hat{P}_{jk} = P_\mu(\hat{X}_{m+1} = k \mid \hat{X}_m = j) \quad (20)$$

Then, following ref 36 theorem 3.1,  $\hat{P}$  can be computed by matrix multiplication using the committors

$$\begin{aligned}\hat{P}_{jk} &= \frac{1}{\hat{\mu}_j} \langle (P^b - \text{Id})q_j^-, q_k^+ \rangle_\mu, \quad j \neq k, \\ \hat{P}_{jj} &= 1 - \sum_{k \neq j} \hat{P}_{jk}\end{aligned}\quad (21)$$

where the inner product is defined by  $\langle u, v \rangle_\mu = \sum_{i \in S} u(i)v(i)\mu(i)$ . In general, the milestone process need not be a Markov process. The results in refs 16 and 34 show, however, that it is an approximate Markov process as long as the core sets are proper metastable sets, that is, if the typical time scale on which  $(\tilde{X}_m)$  leaves  $C$  is much smaller than the typical expected hitting times between the core sets. Thus, by taking  $\hat{P}$  as our MSM transition matrix, we introduce an additional modeling error that is smaller as the core sets become more metastable. With this MSM transition matrix, we can define the MSM kinetics. If we start from some initial probability  $\hat{p}_j(0)$  of being in state  $j$  at time  $t = 0$ , then its evolution  $\hat{p}_j(t)$  in time is discrete in multiples of period  $T$  and given by

$$\hat{p}_j(T) = \sum_k \hat{p}_k(0) \hat{P}_{kj} \quad (22)$$

In our approach  $\hat{p}_j(mT)$  is a good approximation of  $P(\hat{X}_m = j)$  for appropriately chosen core sets.

**Remark 1:** Our definition of the milestoneoning process in terms of the process  $(\tilde{X}_m)$  generated by  $P$  guarantees that we can directly compute  $\hat{P}$  via eq 20 from trajectories of  $(\tilde{X}_m)$  without computing the committor functions. This is of importance if the spatial discretization underlying  $(\tilde{X}_m)$  is fine enough, because the kinetics of  $(\tilde{X}_m)$  then approximates the original kinetics of  $(x_{mT})$ , such that we can directly compute  $\hat{P}$  via eq 20 from NEMD

trajectories without first computing  $P$ , which substantially simplifies MSM building if the core sets are already known.

**Remark 2:** Following ref 16, we can also define another pair of stochastic MSM matrices:

$$\hat{T}_{jk} = \frac{\langle q_j^-, Pq_k^+ \rangle_\mu}{\hat{\mu}_j} \quad (23)$$

$$\hat{M}_{jk} = \frac{\langle q_j^-, q_k^+ \rangle_\mu}{\hat{\mu}_j} \quad (24)$$

that are connected to  $\hat{P}$  by the following identity

$$\hat{P} = \hat{T} - \hat{M} + \text{Id}$$

which can be seen by means of direct computation. For the off-diagonal entries, we have

$$\begin{aligned}\hat{T}_{jk} - \hat{M}_{jk} &= \frac{\langle q_j^-, Pq_k^+ \rangle_\mu}{\hat{\mu}_j} - \frac{\langle q_j^-, q_k^+ \rangle_\mu}{\hat{\mu}_j} \\ &= \frac{\langle (P^b - \text{Id})q_j^-, q_k^+ \rangle_\mu}{\hat{\mu}_j} = \hat{P}_{jk}, \quad j \neq k\end{aligned}\quad (25)$$

Stochasticity yields  $\hat{T}_{jj} - \hat{M}_{jj} + 1 = \hat{P}_{jj}$  for the diagonal entries. Furthermore, as shown in Appendix B,  $\hat{T}$  as well as  $\hat{M}$  can be computed from trajectories without requiring committors.

The importance of the pair  $\hat{T}$  and  $\hat{M}$  for MSM building comes from the following observation: the main NEMD relaxation time scales are given by the dominant eigenvalues of  $P$ ,<sup>16,18</sup> and these dominant eigenvalues can be approximated by discretizing the related eigenvalue problem  $Pu = \lambda u$  by means of a Galerkin approximation with the finite dimensional ansatz space spanned by the forward committors  $q_j^+$ ,  $j = 1, \dots, K$ , together with the finite dimensional test function space spanned by the backward committors  $q_j^-$ ,  $j = 1, \dots, K$  (test functions multiplied from the left by the inner product  $\langle \cdot, \cdot \rangle_\mu$ ). The thus discretized eigenproblem takes the form of a generalized eigenproblem

$$\hat{T}\hat{u} = \hat{\lambda}\hat{M}\hat{u}, \quad \text{or, equivalently} \quad \hat{M}^{-1}\hat{T}\hat{u} = \hat{\lambda}\hat{u}, \quad (26)$$

For the reversible case, it is known that the  $k$  eigenvalues  $\hat{\lambda}$  are very good approximations of the dominant eigenvalues  $\lambda$  of the original problem if the core sets are proper metastable sets.<sup>23</sup> Whether this is true for the nonreversible case is not yet known, but if the deviation from reversibility is weak, and the dominant eigenvalues of  $P$  are real-valued, then the results should hold analogously, see ref 16, theorem 4.19.

If all of its entries are positive such that it is a stochastic matrix,  $\hat{M}^{-1}\hat{T}$  can thus also be taken as MSM transition matrices. In the case of  $\hat{M}^{-1}\hat{T}$ , the MSM modeling error results from Galerkin discretization, whereas the MSM modeling error of  $\hat{P}$  results from ignoring the potential non-Markovianity of  $\hat{X}_m$ .

**Remark 3:** As a matter of fact, if  $\hat{X}_m$  were Markovian, the following identity would hold:  $\hat{P} = \hat{T} \hat{M}^{-1}$  (see Appendix C for the proof). In this case, we would have  $\hat{M}^{-1}\hat{T} = \hat{M}^{-1}\hat{P}\hat{M}$  and the eigenvalues of  $\hat{P}$  and  $\hat{M}^{-1}\hat{T}$  would be identical. Therefore, in practice, the deviation of the eigenvalues of  $\hat{P}$  from those of  $\hat{M}^{-1}\hat{T}$  indicates deviation from Markovianity regarding the process  $\hat{X}_m$ .

#### 4. NUMERICAL EXAMPLE: ALANINE DIPEPTIDE UNDER AN OSCILLATORY ELECTRIC FIELD

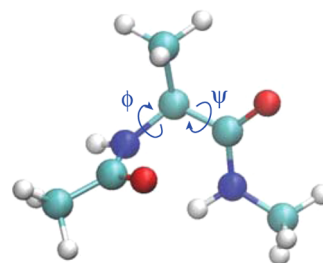
We know that, in theory, whenever the spatial discretization is fine enough, the Markov jump process  $X_t$  associated with Master eq 7 is a good approximation of the original MD process  $x_t$  governed by eq 1. In practice, however, it is difficult to predict how many discrete sets are needed for sufficiently fine discretization. Moreover, because the total dimension of  $x_t$  is  $3N_{\text{atom}}$  ( $N_{\text{atom}}$  being the number of atoms), it is prohibitive to perform an extensive discretization over all degrees of freedom for most systems of practical interest.

One possible way to define appropriate discretization sets is first to find a few collective variables, and then to discretize these collective variables as finely as needed either by uniform or adaptive discretization.<sup>17,37</sup> However, it is difficult to give a general answer a priori regarding how to choose the collective variables and how fine their discretization should be. For large or high dimensional systems, these questions usually become nontrivial.

To illustrate how the discretization works in practice, we take the alanine dipeptide system under an oscillatory EF as an example, the NEMD simulation of which was extensively studied in ref 15. The system was simulated in a  $2.7 \times 2.7 \times 2.7 \text{ nm}^3$  periodic simulation region with one alanine dipeptide molecule, described by the CHARMM27 force field,<sup>38</sup> dissolved in 641 TIP3P water molecules.<sup>39</sup> The grid-based energy correction map (CMAP)<sup>40</sup> was used to correct the backbone dihedral angle energies. All simulations were performed using an in-house modified Gromacs 4.6.5<sup>41</sup> implementing the CHARMM27 force field.<sup>42</sup> The alanine dipeptide was put into the local thermostating environment with a spherical dynamical region of 1.0 nm radius centered at the  $\alpha$  carbon. The Langevin thermostat with target temperature  $\mathcal{T} = 300 \text{ K}$  and time scale  $\tau_T = 0.1 \text{ ps}$  was coupled to the thermostated region. The whole system was coupled to the Parrinello–Rahman barostat<sup>43</sup> (in standard Gromacs implementation) with  $\tau_p = 2.0 \text{ ps}$  to maintain the system at 1 bar. The non-equilibrium trajectories were integrated by the leapfrog scheme with a time step of 0.002 ps. The short-range van der Waals interactions were cutoff at 1.00 nm and were smoothed from 0.95 to 1.00 nm by the “shift” method provided by Gromacs. The energy conserving particle mesh Ewald (PME)<sup>44,45</sup> method (“pme-switch”) was used to compute the long-range electrostatic interaction with the same real-space cutoff radius as the van der Waals interactions. The Gromacs default Fourier spacing of 0.12 nm and B-spline interpolation order of 4 were adopted. The splitting parameter was optimized with respect to the electrostatic force computing accuracy by the Gromacs tool “g\_pme\_error”.<sup>46</sup> The neighbor list was updated every 5 time steps with a list-building radius of 1.20 nm. All hydrogen-involving covalent bonds were constrained by the LINCS algorithm,<sup>47</sup> except for water molecules, which were constrained by the SETTLE algorithm.<sup>48</sup> The whole system was driven by the periodic electric field  $E(t) = E_0 \sin(2\pi t/T)$  and  $D(x) = (1,0,0)^T$  with the intensity of the field being  $E_0 = 1.0 \text{ V/nm}$  and the period being  $T = 10 \text{ ps}$ . The 20,000 branching trajectories were simulated from 20,000 initial configurations that sampled the equilibrium distribution. The equilibrium configurations were prepared by an equilibrium MD simulation of  $10^6 \text{ ps}$  in length, along which snapshots were saved every 50 ps. The branching NEMD trajectories were each 4,000 ps long, and the system reached a non-equilibrium quasi-stationary state in roughly 300 ps.

For this periodically driven molecular system, we will first show how to choose an appropriately fine spatial discretization. After validating this discretization, we will consider the time-discretized dynamics generated by the Floquet transition matrix  $P$  in comparison to the original NEMD simulation. Finally, we will coarse-grain this description further by construction of a 3-state Markov state model that is able to describe the long-term kinetics of the system correctly.

**4.1. Spatial Discretization.** We choose the dihedral angles  $\phi$  and  $\psi$  as collective variables (see Figure 1), and the



**Figure 1.** A schematic plot of the alanine dipeptide molecule and the dihedral angles  $\phi$  and  $\psi$ .

discretization is a uniform partition of the  $\phi$ – $\psi$  plane. We denote the number of discretization intervals on each dihedral by  $N_{\text{dih}}$ , yielding  $N = N_{\text{dih}}^2$  discretization sets  $\{\Omega_i\}$ ,  $i \in S = \{1, \dots, N_{\text{dih}}^2\}$ .

On the basis of a given spatial discretization, we can aggregate the transition matrix  $P = \Phi^T(T)$  simply by counting the transition behavior of MD trajectories

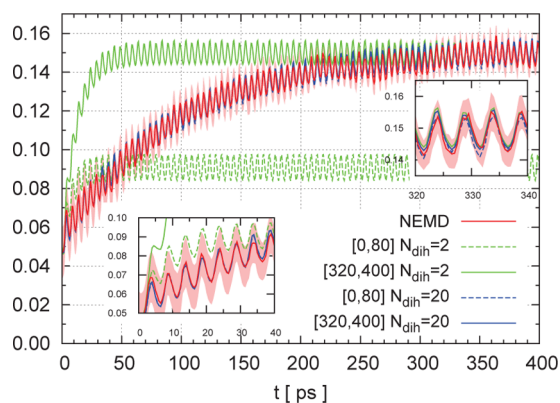
$$P_{ij} = P(X_T = j | X_0 = i) \quad (27)$$

However,  $P$  allows approximation of the original dynamics on multiple  $mT$  of the period only. To have a time-continuous description, we need the generator  $L(t)$  of the Master equation. If the discretization is fine enough, one possible approximation to  $L(t)$  is via the following forward finite difference scheme

$$L_{ij}(t) \approx \frac{1}{\tau} [P(X_{t+\tau} = j | X_t = i) - \delta_{ji}], \quad i, j \in S \quad (28)$$

with  $\tau$  as an appropriate small enough lag-time. Because the dimensionality is reduced by using only a few collective variables, the lag-time should be chosen sufficiently large that the original dynamics  $x_t$  is properly relaxed with regard to the unresolved degrees of freedoms on time scales shorter than the lag-time (assuming that the collective variables capture the slow dynamics). Below, we investigate the discretization quality with respect to the choice of  $N_{\text{dih}}$  and lag-time  $\tau$ .

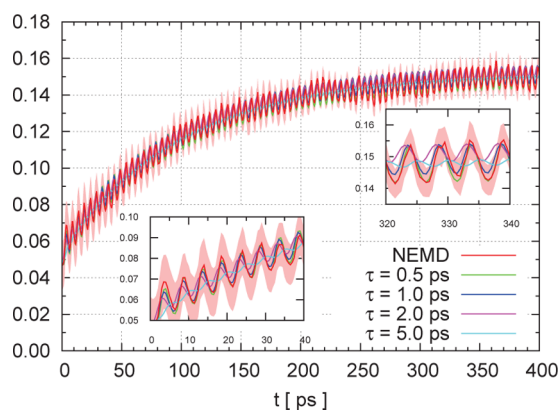
We estimate the discretized generator  $L(t)$ ,  $t \in [0, T)$  from NEMD trajectories generated by  $\mathcal{L}(t)$  in different time intervals  $[t_1, t_2]$ . As discussed above, whenever the discretized dynamics approximates the original dynamics well, the time-periodic generator  $L(t)$  should not depend on the choice of the interval  $[t_1, t_2]$  in the estimation procedure in sch 28, provided that the initial state of the system is not very far from the stationary state at a long-time limit. Therefore, this is an indicator for calibrating the discretization quality. We compute  $L(t)$  by two discretizations,  $N_{\text{dih}} = 2$  and  $N_{\text{dih}} = 20$ , and two choices of time intervals,  $[0, 80] \text{ ps}$  and  $[320, 400] \text{ ps}$ , and then compare the time-dependent probability  $P(\phi_i \in [0, 180), \psi_i \in [0, 180))$  with the brute force NEMD result in Figure 2 using a lag-time of  $\tau = 0.5 \text{ ps}$ . Using  $N_{\text{dih}} = 2$ , the dynamics depends on the time interval used for



**Figure 2.** Time-dependent probability  $P(\phi_t \in [0,180], \psi_t \in [0,180])$ . The brute force NEMD simulation is compared with different spatial discretization methods. The red shadow region indicates the statistical uncertainty of the NEMD simulation.

calculating the generator. Using a time interval of  $[0, 80]$  ps, the discretized dynamics deviates from the NEMD result, whereas using the time interval  $[320, 400]$  ps the discretized dynamics can only reproduce the NEMD result after 300 ps. This therefore indicates poor approximations of the original dynamics with  $N_{\text{dih}} = 2$ . The reason for this is that the discretization with  $N_{\text{dih}} = 2$  is too coarse such that the dynamics cannot be fully equilibrated within the lag-time  $\tau$  in each discretized set; therefore, the discretization presents state dependency. For  $N_{\text{dih}} = 20$ , the discretized dynamics does not depend on the time interval of calculating the generator and is consistent with the NEMD simulation within the error bars. Therefore, throughout this paper, we use  $N_{\text{dih}} = 20$  to discretize the dihedral angle space of the alanine dipeptide.

Next, we discuss the effect of lag time  $\tau$  on the estimation of the generator. Therefore, we consider different choices of  $\tau$  (0.5, 1.0, 2.0, and 5.0 ps; see Figure 3), all based on the identical dihedral



**Figure 3.** Time-dependent probability  $P(\phi \in [0,180], \psi \in [0,180])$ . The NEMD simulation is compared with the Master equation using generators discretized with  $N_{\text{dih}} = 20$  and different lag times ( $\tau$ ). The red shadow region indicates the statistical uncertainty of the NEMD simulation.

angle discretization using  $N_{\text{dih}} = 20$ . It is clear that when the lag-time is close to the period (10 ps) the discretized dynamics cannot resolve the probability change within a period. However, it is surprising that even quite large lag-times are capable of capturing the overall long-term behavior of the original dynamics. We observe no significant difference between

$\tau = 0.5$  and  $\tau = 1.0$  ps, which means that the discretized dynamics is not very sensitive to the choice of  $\tau$ . Therefore,  $\tau = 0.5$  ps will be used throughout this paper.

**4.2. Quasi-Stationary Distribution ( $\mu$ ).** After having validated the fine-scale spatial discretization, we will now consider the time-homogeneous process  $\tilde{X}_m$  generated by the Floquet transition matrix  $P = \Phi^T(T)$  and investigate whether it reproduces the properties of the original non-equilibrium process  $x_t$ . In this context, only the configurations at the integral periods  $mT$  along the original process are taken into consideration.

An important check is the consistency between the stationary probability density of  $\Phi(T)$  (i.e., the leading eigenvector  $\mu$ ) and that estimated from the original NEMD simulation

$$\rho_{\text{st}}(\phi, \psi) = \lim_{m \rightarrow \infty} \rho(\phi, \psi, mT) \quad (29)$$

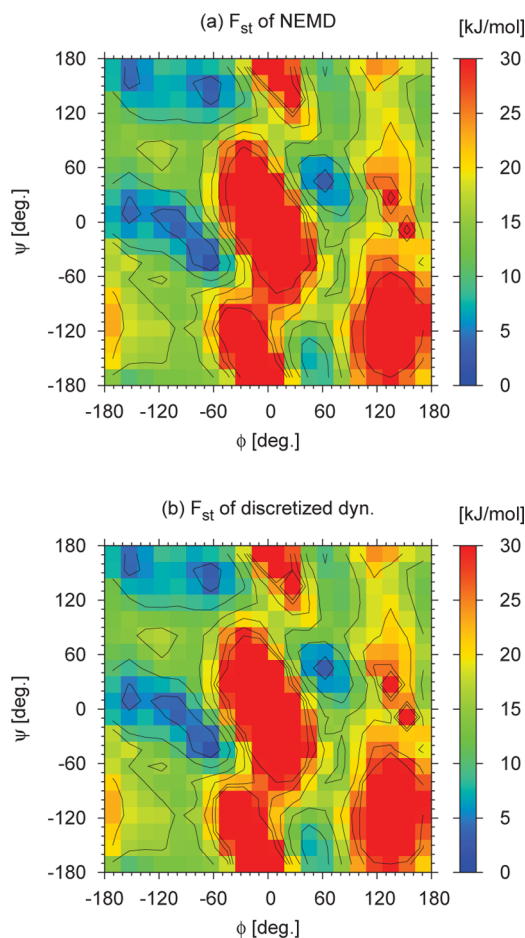
On each NEMD branching trajectory, the initial 320 ps are discarded, and the rest of the trajectory in the time interval  $[320, 4000]$  ps is averaged to estimate the quasi-stationary probability distribution  $\rho_{\text{st}}$ .  $P$  is computed as described above, and then  $\mu$  is computed as its leading eigenvector. To make it comparable to the free energy in the equilibrium case, we take the logarithm of the distributions, namely,  $F_{\text{st}}(\phi, \psi) = -k_B T \log \rho_{\text{st}}(\phi, \psi)$  for NEMD and  $F_{\text{st}}(\phi, \psi) = -k_B T \log \mu(\phi, \psi)$  for  $P$ , where  $k_B$  is the Boltzmann constant and  $T$  is the temperature of the system. The results are compared in Figure 4, where good consistency between the NEMD simulation and  $P$  is observed.

**4.3. Core Set Identification.** The procedure for identifying good metastable core sets of the irreversible Markov process associated with  $P$  is described in detail in ref 34; here, we provide only the fundamental idea behind it. If strong metastable sets  $C_j$ ,  $j = 1, \dots, K$  exist, they should have one main property: when starting from a state in  $C_j$ , the expected hitting time of a state in  $C_i$  should be much shorter than that of any state in one of the other sets  $C_j$ ,  $j \neq i$ . In fact, the hitting time distribution should exhibit roughly constant levels in each set  $C_j$  and should vary significantly in the transition region  $C = \Omega \setminus \bigcup_j C_j$  between the metastable sets. If starting from some randomly chosen initial states, one can thus identify the metastable core sets and the transition region by analyzing the hitting time distributions. This procedure is similar to the procedures used for reversible processes<sup>17,18,49</sup> but utilizes hitting time distributions instead of any eigenvector information.

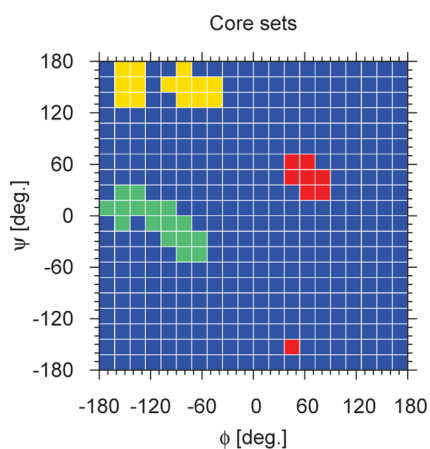
The metastable core sets identified by this procedure based on the estimation of  $P$  are illustrated in Figure 5 and are denoted by  $C_{\alpha_R}$  (green),  $C_\beta$  (yellow), and  $C_{\alpha_L}$  (red). They correspond to the centers of the wells in the free energy landscapes shown in Figure 4 and to the right-handed  $\alpha$  helix,  $\beta$  sheet, and left-handed  $\alpha$  helix conformations of the peptide, respectively.

**4.4. First Mean Hitting Times.** The first mean hitting time as a function of the dihedral angles  $(\phi, \psi)$ , is defined by the expected first time needed for hitting a certain core set  $C_j$ ,  $j \in \{\alpha_R, \beta, \alpha_L\}$  conditioned on starting from the conformation  $(\phi, \psi)$ , or more exactly, from equilibrium conformations  $(\phi, \psi) \in \Omega_i$ . Because the largest first mean hitting time (starting from states in core set  $\alpha_R$  and hitting  $\alpha_L$ ) is longer than 600 ps, the results will be biased if we use the NEMD trajectories of 4000 ps in length for brute force Monte Carlo estimation of the hitting time. Therefore, we base our Monte Carlo estimate on 100 NEMD trajectories of  $2 \times 10^5$  ps instead. For comparison, we compute the first mean hitting time of the discretized dynamics via its transition matrix  $P$ . The first mean hitting time  $h_{C_j}(i)$  of core set





**Figure 4.** Color-scale plot of the logarithmic quasi-stationary distribution  $\mu$  of the (a) NEMD and (b) discretized dynamics governed by the Floquet transition matrix  $P = \Phi^T(T)$ .

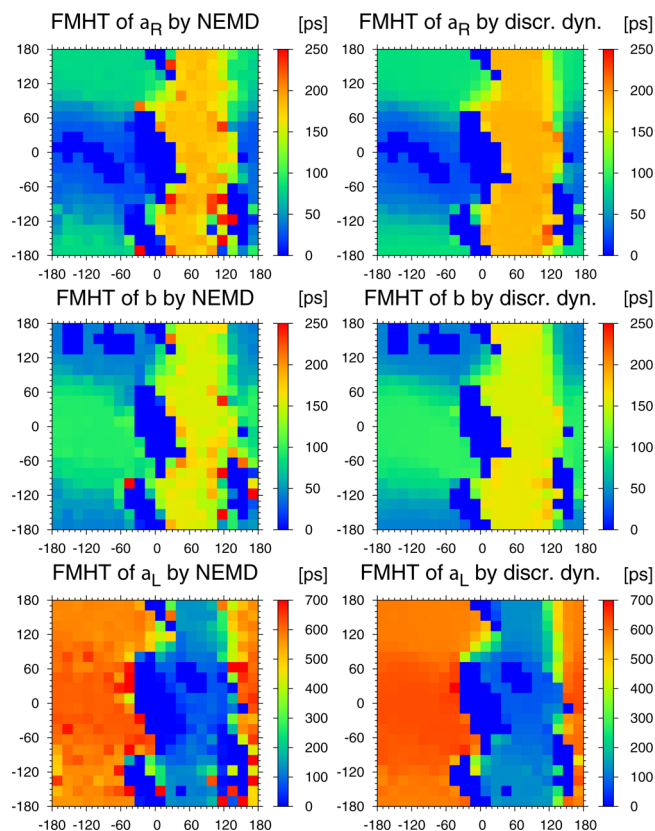


**Figure 5.** Core set identification. Different colors indicate different core sets:  $C_{\alpha_R}$  (green),  $C_\beta$  (yellow), and  $C_{\alpha_L}$  (red). The blue color indicates the transition region ( $C$ ) that does not belong to one of the core sets.

$C_j$  starting in  $\Omega_i$  can be computed by means of solving the linear problem<sup>16</sup>

$$(P - \text{Id})h_{C_j}(i) = -1, \quad \text{if } C_j \cap \Omega_i = \emptyset$$

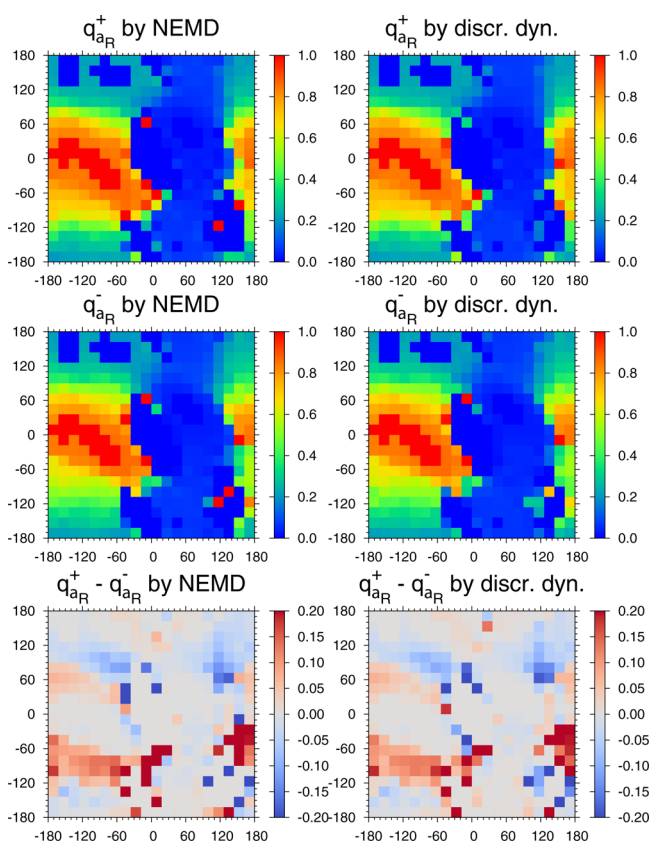
The resulting first mean hitting times are presented in Figure 6. The good consistency between the NEMD estimate and the discretized Markov process  $\tilde{X}_m$  indicates a good approximation.



**Figure 6.** Comparisons of the first mean hitting time (FMHT) based on an NEMD simulation (left column) and discretized dynamics (right column). From top to bottom, the first mean hitting times to the core sets  $C_{\alpha_R}$ ,  $C_\beta$  and  $C_{\alpha_L}$  are shown, respectively.

One should note that the NEMD estimate of the first mean hitting time is subject to statistical sampling errors, whereas the first mean hitting times  $h_{C_j}$  only contain statistical errors coming from the estimation of  $P$ . Thus, using  $P$  helps in calculating the observables more smoothly (i.e., less statistical error, no additional sampling). Additionally, the computational cost of the discretized process, if the cost for estimating  $\Phi(T)$  is not included, is essentially smaller than NEMD. The computation of  $h_{C_j}$  was a matter of milliseconds on a laptop, whereas the NEMD trajectories took  $1.6 \times 10^4$  core hours on Intel Xeon E5-4650 CPUs.

**4.5. Forward and Backward Committors.** Committors are very important statistical properties of Markov processes<sup>24,50</sup> and play an important role in MSM building<sup>16,29,36</sup> (see below). Therefore, it is worth checking if the discretized process  $\tilde{X}_m$  reproduces the NEMD committors. The forward committor  $q_j^+(i)$  of a core set  $C_j$ ,  $j \in \{\alpha_R, \beta, \alpha_L\}$  is defined as the probability of visiting core set  $C_j$  next conditioned on starting at conformation  $(\phi, \psi) \in \Omega_i$ . The backward committor  $q_j^-(i)$  of a core set  $C_j$ ,  $j \in \{\alpha_R, \beta, \alpha_L\}$  is defined as the probability of last coming from  $C_j$  conditioned on having arrived at configuration  $(\phi, \psi) \in \Omega_i$  at this time. For reversible Markov processes, the forward and backward committors are identical; however, this is generally not the case for irreversible processes. The committors estimated from NEMD simulations (20,000 trajectories, 4,000 ps each) are compared with those computed from  $P$  by means of solving the linear eqs 18 and 19. Figures 7–9 present both committors as well as their differences corresponding to the different core sets. The committors of the discretized process have good consistency with those of the NEMD simulations. The non-zero values in the

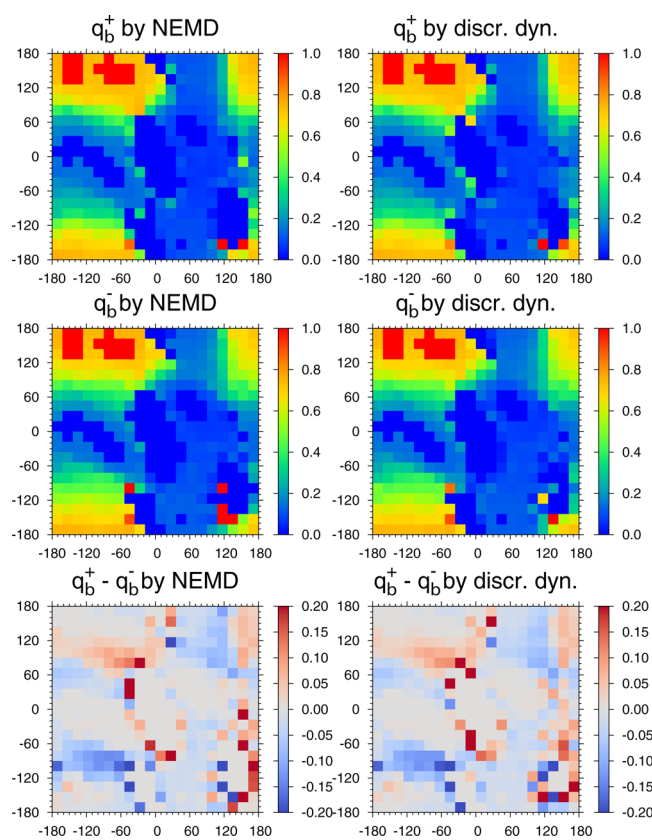


**Figure 7.** Forward  $q_{a_R}^+$  and backward  $q_{a_R}^-$  committors and their difference  $q_{a_R}^+ - q_{a_R}^-$  computed from the NEMD trajectories (left column) and the discretized dynamics (right column).

committor differences indicate that the NEMD process, projected on the discretized dihedral angle space, is irreversible, and that the discretized process is able to correctly describe this reversibility. In addition, and subsequently of central importance, the accurate reproduction of the committors indicates that it is reasonable to build the MSM out of the committors of the discretized process.

**4.6. MSM Building and Validation.** Following the process described in section 3, we are able to build a 3-state MSM for the externally driven alanine dipeptide system, where the quasi-stationary probability distribution  $\mu$ , the three core sets, and the forward and backward committors are estimated as described above. The MSM transition matrices  $\hat{P}$ ,  $\hat{M}$ , and  $\hat{T}$  are then evaluated using eqs 21, 23, and 24, respectively. Alternatively, the MSM transition matrix  $\hat{P}$  is calculated directly from the NEMD trajectories using eq 20 (very good agreement with the one computed from the committors). The leading eigenvalues of  $P = \Phi^T(T)$  are compared with those of  $\hat{P}$  and  $\hat{M}^{-1}\hat{T}$  in Table 1. Not surprisingly, the two approaches for MSM building are consistent. The MSM is able to accurately reproduce the largest nontrivial eigenvalue, which means a precise reproduction of the longest nontrivial implied time scale. The accuracy of the second nontrivial time scale is not as good as the first but is still acceptable.

The lower accuracy of the second nontrivial time scale stems from the time-discretization underlying our approach. The associated problem is graphically illustrated in Figure 10. The milestone process is defined to change state when the trajectory hits a core set that is not the one it came from. In Figure 10, the milestone process based on the time-continuous trajectory



**Figure 8.** Forward  $q_b^+$  and backward  $q_b^-$  committors and their difference  $q_b^+ - q_b^-$  computed from NEMD trajectories (left column) and the discretized dynamics (right column).

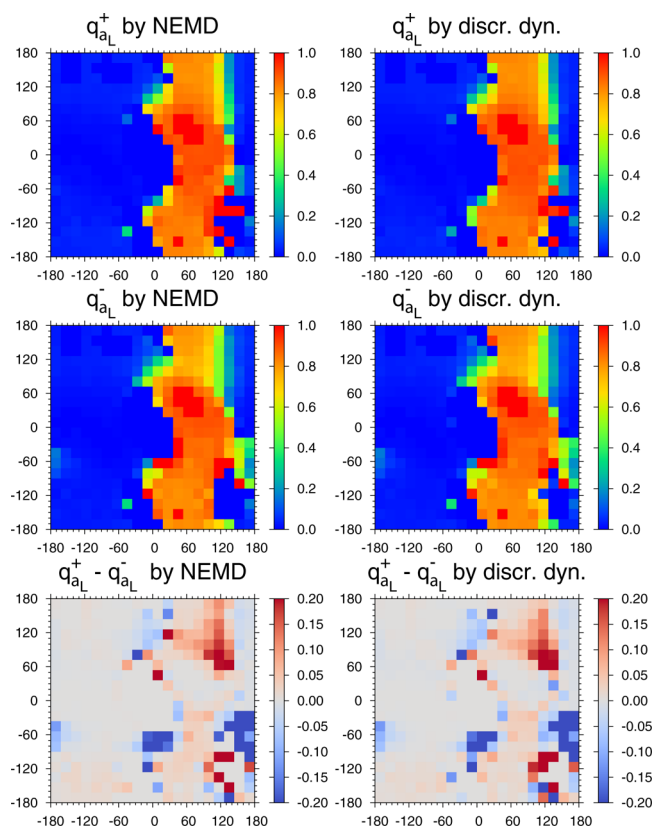
changes to state “2” at the first hit of core set  $C_2$ . In comparison, this first hitting event would not be correctly recognized using the time-discretized trajectory; instead, in the time-discretized case, the milestone process would change to state “2” only at the second hit of core set  $C_2$ . This respective error in detecting the correct first hitting time will not be significant if it is much longer than the discretization time step. However, it will become relevant if both are comparable. In our case, the temporal discretization is given by the period  $T = 10$  ps, which is not well separated from the second nontrivial time-scale that is 26.5 ps (calculated by  $-T/\log(\lambda_2)$ ). Therefore, we observe an unavoidable deviation in the second time scale. In applications of more complex and realistic problems (especially larger systems), particularly for systems exhibiting time scale separation between the external driving period and the interested implied time scales, good accuracy is expected.

The difference between the eigenvalues of  $\hat{P}$  and  $\hat{M}^{-1}\hat{T}$  can be taken as an indication of the non-Markovianity of the milestone process  $\hat{X}_m$ . This non-Markovianity seems to have a stronger influence on the second nontrivial time scale than on the first; this may be caused by shorter decorrelation time scales due to weaker metastability of the core sets involved. Numerically, the two MSM transition matrices are

$$\hat{M}^{-1}\hat{T} = \begin{pmatrix} 0.860 & 0.133 & 0.008 \\ 0.192 & 0.775 & 0.033 \\ 0.019 & 0.066 & 0.915 \end{pmatrix}$$

$$\hat{P} = \begin{pmatrix} 0.882 & 0.110 & 0.008 \\ 0.158 & 0.815 & 0.028 \\ 0.023 & 0.055 & 0.922 \end{pmatrix}$$

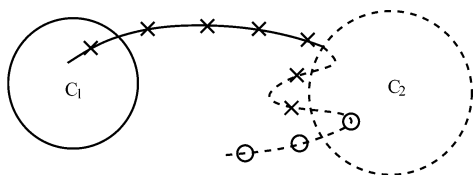




**Figure 9.** Forward  $q_{a_L}^+$  and backward  $q_{a_L}^-$  committors and their difference  $q_{a_L}^+ - q_{a_L}^-$  computed from NEMD trajectories (left column) and the discretized dynamics (right column).

**Table 1.** Comparison of Second and Third Eigenvalues of  $P$  and  $3 \times 3$  MSM Transition Matrices  $\hat{P}$  and  $\hat{M}^{-1}\hat{T}$ , Respectively, from the Two MSM Approaches

	$\lambda_2$	$\lambda_3$	$\lambda_4$
$P$	0.905	0.668	0.551
$\hat{P}$	0.909	0.710	
$\hat{M}^{-1}\hat{T}$	0.901	0.649	



**Figure 10.** Schematic plot of the comparison between continuous and time-discretized trajectories in hitting a core set. The time-continuous trajectory changes to state “2” at the first hit of the core set  $C_2$  (indicated by changing from solid to dashed line). The time-discretized trajectory changes to state “2” at the second hit (indicated by changing from “x” to “o”).

from which we see that the left-handed  $\alpha$  helix conformations of the peptide exhibit the strongest metastability.

It is worth noting that, although the discretized process  $\tilde{X}_m$  is irreversible, the MSM built out of it is almost reversible. The magnitude of the antisymmetric part of the matrix  $\text{diag}(\hat{\mu})\hat{P}$  is only on the order of  $10^{-4}$ .

In fact,  $\hat{P}$  can be considered the fingerprint of the long-term kinetics (see refs 26 and 27) of the alanine dipeptide in an

oscillatory electric field. To provide further validation of this statement, we study time-dependent expectation values of the form

$$\mathcal{A}(t) = \langle A(i) \rangle_t = \sum_{i \in S} A(i) p(i, t) \quad (30)$$

where  $p(i, t) = P(X_t = i)$  is the probability of being in set  $\Omega_i$  at time  $t$  as governed by the Master equation, and the observable  $\mathcal{A}$  is spanned by the backward committors, namely,

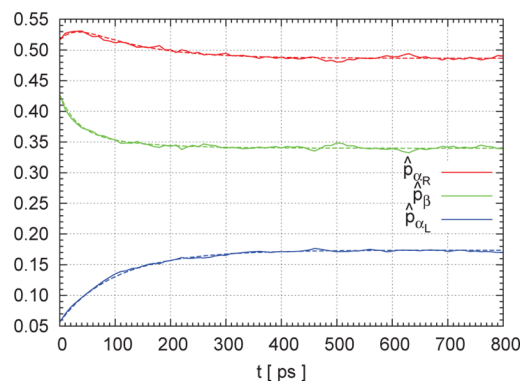
$$A(i) = \sum_{j=1}^K \alpha_j q_j^-(i) \quad (31)$$

Then

$$\begin{aligned} \mathcal{A}(t) &= \sum_{i \in S} \sum_{j=1}^K \alpha_j q_j^-(i) p(i, t) \\ &= \sum_{i \in S} \sum_{j=1}^K \alpha_j P(\hat{X}_t = j | X_t = i) P(X_t = i) \\ &= \sum_{i \in S} \sum_{j=1}^K \alpha_j P(\hat{X}_t = j, X_t = i) \\ &= \sum_{j=1}^K \alpha_j P(\hat{X}_t = j) = \sum_{j=1}^K \alpha_j \hat{p}_j(t) \end{aligned} \quad (32)$$

where the time-dependent probability  $\hat{p}_j(t)$  of being assigned to MSM macrostate  $j$  at time  $t$  can be computed by means of the MSM via simple matrix multiplications using eq 22.

In Figure 11, we compare the numerical calculation of  $\hat{p}_j(mT)$ ,  $m \in \mathbb{N}$  from NEMD and MSM calculations. In the



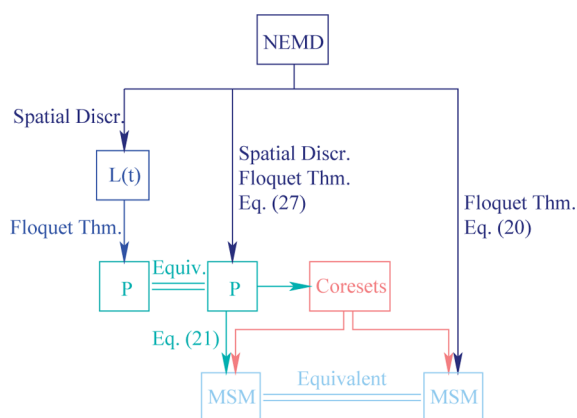
**Figure 11.** Comparison of NEMD and MSM long-term kinetics of alanine dipeptide in an oscillatory electric field. The plots show the time-dependent probability  $\hat{p}_j(t)$  to be assigned to core set  $C_j$  (corresponding to the observable  $\mathcal{A}$  given in eq 31 with  $\alpha_j = 1$  and  $\alpha_k = 0$  for  $k \neq j$ ). Solid lines are from brute force NEMD simulations, whereas the dashed lines are from our MSM.

NEMD case, the identity  $\hat{p}_j(mT) = \sum_{i \in S} q_j^-(i) p(i, mT)$  is used, and the backward committor and the probability density on the right hand side are estimated directly from the NEMD trajectories. For the MSM, projection of the initial probability is applied,  $\hat{p}_j(0) = \sum_{i \in S} q_j^-(i) p(i, 0)$ , then the time-dependent probability at  $mT$  is generated by eq 22 (i.e., by simple matrix multiplication). The agreement is almost perfect.

## 5. CONCLUDING REMARKS AND DISCUSSION

In this paper, we proposed methods for MSM building for a periodically driven non-equilibrium system. We demonstrated their validity and performance by application to alanine dipeptide under an oscillatory electric field. We showed that the proposed methods allow for capture of the long-term behavior of the original non-equilibrium dynamics.

We provided effective methods for discretizing the original NEMD dynamics both temporally and spatially; the end-product is a time-homogeneous, and generally irreversible, Markov jump process. Discretization was done via two equivalent approaches: either by a two-step approach with first spatial and then temporal discretization or a one-step approach that involves both discretizations simultaneously. These two version are shown by the left-most and middle branches of the diagram in Figure 12.



**Figure 12.** Flowchart showing optional procedures for MSM building based on NEMD trajectories. A lighter color indicates a “coarser” approximation of the original NEMD dynamics. Please note that “equivalence” only means that the respective procedures result in the same matrix/MSM if one assumes perfect sampling.

Although the end-product of the two-step and one-step discretization procedures are formally equivalent, it is clear that the one-step discretization does not preserve any information within one period because only the states at integral multiples of the period of the external forcing are considered. This is not a serious problem whenever the long-term behavior of the system is of interest and the corresponding time scales are significantly longer than one period. However, if the time scale of interest is comparable to the period, the two-step discretization is preferable because it allows recovery of the dynamics between multiples of the period.

Building the final MSM based on the time-homogeneous discretized dynamics is straightforward by using eq 21 and a set of core sets that are derived from the discretized dynamics. In application to the alanine dipeptide, numerical results show that a 3-state MSM can reproduce the leading nontrivial eigenvalue with very good accuracy and the second nontrivial eigenvalue with acceptable accuracy. The lower accuracy of the second nontrivial eigenvalue may result from the fact that the second slowest time scale is not significantly longer than the period. By means of this 3-state MSM, we can reproduce with almost perfect accuracy the time-evolution of the population of the main conformations induced by the periodic forcing when starting from the equilibrium distribution of the unforced molecular system.

The right-most branch in Figure 12 presents an equivalent, and seemingly much simpler alternative to the middle branch: MSM building directly from NEMD trajectories. In practice, however, this method may not be applicable because it requires predefined core sets, and the identification of core sets is usually not a trivial task, especially for molecular dynamics under non-equilibrium conditions. This task is substantially simplified when an accurate time-homogeneous discretization of the original NEMD process is available (that is, has been constructed in advance). In the present study, we computed the core sets by finding almost constant levels of the hitting time distribution for the discretized dynamics. This procedure is itself a novelty because it does not require any spectral information, such as eigenvectors, as standard approaches do (see refs 17 and 18).

In this paper, we mainly focus on the development of the first available methods for MSM building in non-equilibrium systems. However, the application of these methods to the conformation dynamics of the alanine dipeptide results in some observations that are interesting in themselves. Under an oscillatory EF, the population of the left-handed  $\alpha$ -helical conformation significantly increases relative to the equilibrium population (see discussions in ref 15), and the leading relaxation time scale of the system is much shorter than in the equilibrium case.

A final remark concerning the utilization of the proposed methodology may be in order. The approach presented herein allows for MSM building for periodic forcings only; an extension to nonperiodic driving forces is not straightforward because in general Floquet theory is not applicable to nonperiodic driving forces. However, from an MSM for a given periodic forcing, optimal control problems may come into play. Using available methods, such as non-equilibrium linear response theory<sup>51</sup> or computational alchemy for MSMs,<sup>52</sup> one can construct the MSM for slightly changed parameters (e.g., period and amplitude) of the external forcing by appropriate reweighting of the MSM for the given forcing. This, in principle, allows for answering questions such as the following: for which parameters of the external forcing does one achieve maximal population of the left-handed  $\alpha$ -helix? Such questions, however, will be addressed the future studies.

## APPENDIX A. REVERSIBILITY OF THE ORIGINAL DYNAMICS

We consider the governing dynamics eq 1. For simplicity, we denote the force by  $F(x_t, t) = -\nabla_x V(x_t) + E(t)D(x_t)$ . We denote  $\sigma = (2\beta^{-1})^{1/2}$ . According to Girsanov, we have

$$\frac{dp[x_t]}{dw[x_t]} = \exp\left\{\frac{1}{\sigma^2} \int_0^T F(x_t, t) dx_t - \frac{1}{2\sigma^2} \int_0^T F^2(x_t, t) dt\right\} \quad (33)$$

where  $dp$  is the probability measure of trajectory  $x_t$  and  $dw$  is the probability measure of the standard Wiener process  $dx_t = \sigma dw_t$ . Assuming discretization of the stochastic process at time  $0 < t_1 < t_2 < \dots < t_N = T$ , where  $t_i = iT/N$ . We denote  $x_i = x_{t_i}$  and  $w_i = w_{t_i}$ , then in the sense of Ito, we have

$$\frac{dp[x_t]}{dw[x_t]} \approx \exp\left\{\frac{1}{\sigma^2} \sum_{i=0}^{N-1} F(x_i, t_i)(x_{i+1} - x_i) - \frac{1}{2\sigma^2} \sum_{i=0}^{N-1} F^2(x_i, t_i)\Delta t\right\} \quad (34)$$

Now, consider a conjugate trajectory  $x_t^\dagger = x_{T-t}$  that starts at  $x_T$  and ends at  $x_0$ . The conjugate dynamics is driven by

$F^\dagger(x_i^\dagger, t) = F(x_i^\dagger, T-t)$ . Writing the Girsanov for the conjugate dynamics

$$\begin{aligned} \frac{dp^\dagger[x_t^\dagger]}{dw[x_t^\dagger]} &\approx \exp\left\{\frac{1}{\sigma^2} \sum_{i=0}^{N-1} F^\dagger(x_i^\dagger, t_i)(x_{i+1}^\dagger - x_i^\dagger) \right. \\ &\quad \left. - \frac{1}{2\sigma^2} \sum_{i=0}^{N-1} [F^\dagger(x_i^\dagger, t_i)]^2 \Delta t\right\} \\ &= \exp\left\{\frac{1}{\sigma^2} \sum_{i=0}^{N-1} F(x_i^\dagger, T-t_i)(x_{i+1}^\dagger - x_i^\dagger) \right. \\ &\quad \left. - \frac{1}{2\sigma^2} \sum_{i=0}^{N-1} [F(x_i^\dagger, T-t_i)]^2 \Delta t\right\} \\ &= \exp\left\{\frac{1}{\sigma^2} \sum_{i=0}^{N-1} F(x_{N-i}, t_{N-i})(x_{N-i-1} - x_{N-i}) \right. \\ &\quad \left. - \frac{1}{2\sigma^2} \sum_{i=0}^{N-1} [F(x_{N-i}, t_{N-i})]^2 \Delta t\right\} \end{aligned} \quad (35)$$

$$\begin{aligned} \frac{dp^\dagger[x_t^\dagger]}{dw[x_t^\dagger]} &= \exp\left\{\frac{1}{\sigma^2} \sum_{i=N}^1 F(x_i, t_i)(x_{i-1} - x_i) \right. \\ &\quad \left. - \frac{1}{2\sigma^2} \sum_{i=N}^1 F^2(x_i, t_i) \Delta t\right\} \end{aligned} \quad (36)$$

Because it is obvious that  $dw[x_t^\dagger]/dw[x_t] = 1$ ,

$$\frac{dp^\dagger[x_t^\dagger]}{dw[x_t^\dagger]} \approx \exp\left\{\frac{1}{\sigma^2} \sum_{i=1}^N F(x_i, t_i)(x_{i-1} - x_i) - \frac{1}{2\sigma^2} \sum_{i=1}^N F^2(x_i, t_i) \Delta t\right\} \quad (37)$$

The difference between the single trajectory probabilities is

$$\frac{dp^\dagger[x_t^\dagger]}{dp[x_t]} \approx \exp\left\{-\frac{1}{\sigma^2} \sum_{i=1}^{N-1} [F(x_i, t_i)(x_{i+1} - x_i) + F(x_i, t_i)(x_i - x_{i-1})]\right\} \quad (38)$$

Assuming the smoothness of the external perturbation, consider the differentiation

$$\begin{aligned} F(x_i, t_i) - F(x_{i-1}, t_{i-1}) &= F(x_i, t_i) - F(x_{i-1}, t_i) + F(x_{i-1}, t_i) - F(x_{i-1}, t_{i-1}) \\ &= \nabla_x F(x_{i-1}, t_i)(x_i - x_{i-1}) + O(\Delta t) \\ &= \nabla_x F(x_{i-1}, t_{i-1})(x_i - x_{i-1}) + O(\Delta t) \end{aligned} \quad (39)$$

The second order expansion with respect to  $x_i - x_{i-1}$  is of the order  $\Delta t$ , so it is absorbed into  $O(\Delta t)$ . Then eq 38 becomes

$$\begin{aligned} \frac{dp^\dagger[x_t^\dagger]}{dp[x_t]} &\approx \exp\left\{-\frac{2}{\sigma^2} \sum_{i=0}^{N-1} F(x_i, t_i)(x_{i+1} - x_i) \right. \\ &\quad \left. - \frac{1}{\sigma^2} \sum_{i=0}^{N-1} \nabla_x F(x_i, t_i)(x_{i+1} - x_i)^2\right\} \end{aligned} \quad (40)$$

Using the identity  $dt = (dw_t)^2 = \sigma^{-2} dx_t^2$ , eq 40 is written in the integral form

$$\frac{dp^\dagger[x_t^\dagger]}{dp[x_t]} \approx \exp\left\{-\frac{2}{\sigma^2} \int_0^T F(x_t, t) dx_t - \int_0^T \nabla_x F(x_t, t) dt\right\} \quad (41)$$

One would not have the second integral on the exponent if the first integral of the exponent were defined in the sense of Stratonovich.

We notice that

$$\begin{aligned} dV(x, t) &= \frac{\partial V}{\partial x} dx + \frac{\partial V}{\partial t} dt \\ &= \frac{1}{2} \sigma^2 \nabla_x^2 V dt + \nabla V dx_t + \frac{\partial V}{\partial t} dt \\ &= -\frac{1}{2} \sigma^2 \nabla_x F dt - F dx_t + \frac{\partial V}{\partial t} dt \end{aligned} \quad (42)$$

Eq 41 becomes

$$\frac{dp^\dagger[x_t^\dagger]}{dp[x_t]} = \exp\left\{\frac{2}{\sigma^2} \left[V(x_T, T) - V(x_0, t_0) - \int_0^T \partial_t V(x_t, t) dt\right]\right\} \quad (43)$$

Taking the limit of the infinite small time interval, noticing that the equilibrium invariant probability density with respect to potential  $V(x,0)$  satisfies  $\mu(x) \propto e^{-\beta V(x,0)}$ , and replacing  $\sigma^2$  by  $2\beta^{-1}$  results in

$$\frac{dp^\dagger[x_t^\dagger]}{dp[x_t]} = \frac{\mu(x_0)}{\mu(x_T)} \exp\left\{-\beta \int_0^T \partial_t V(x_t, t) dt\right\} \quad (44)$$

### A.1. Irreversibility of the Periodic Time-Symmetric Dynamics

In eq 35, we assume periodicity of the perturbation  $F(x, t) = F(x, t+T)$  and symmetry of the external perturbation (i.e.,  $F(x, -t) = F(x, t)$ ), giving

$$\begin{aligned} \frac{dp^\dagger[x_t^\dagger]}{dw[x_t^\dagger]} &\approx \exp\left\{\frac{1}{\sigma^2} \sum_{i=0}^{N-1} F(x_i^\dagger, t_i)(x_{i+1}^\dagger - x_i^\dagger) \right. \\ &\quad \left. - \frac{1}{2\sigma^2} \sum_{i=0}^{N-1} [F(x_i^\dagger, t_i)]^2 \Delta t\right\} \end{aligned} \quad (45)$$

By changing notation  $x^\dagger$  back to  $x$  and comparing with 34, the reversed dynamics is subject to eq 1 (i.e.,  $dp^\dagger = dp$ ). Therefore

$$\begin{aligned} p(x_0, T|x_T, 0) &= \int_{C\{x_T, 0; x_0, T\}} dp[x_t^\dagger] \\ &= \int_{C\{x_0, 0; x_T, T\}} \frac{dp[x_t^\dagger]}{dp[x_t]} dp[x_t] \\ &= \int_{C\{x_0, 0; x_T, T\}} \frac{dp^\dagger[x_t^\dagger]}{dp[x_t]} dp[x_t] \\ &= \frac{\mu(x_0)}{\mu(x_T)} \int_{C\{x_0, 0; x_T, T\}} \\ &\quad \times \exp\left\{-\beta \int_0^T \partial_t V(x_t, t) dt\right\} dp[x_t] \end{aligned} \quad (46)$$



where  $C\{x_0, 0; x_T, T\}$  denotes all continuous trajectories starting at  $x_0$  and ending at  $x_T$ . If  $\partial_t V = 0$  (i.e., equilibrium), we get

$$p(x_0, T|x_T, 0)e^{-\beta V(x_T, T)} = p(x_T, T|x_0, 0)e^{-\beta V(x_0, 0)} \quad (47)$$

which proves the reversibility of the equilibrium dynamics. The term

$$W[x_t] = \int_0^T \partial_t V(x_t, t) dt \quad (48)$$

is the non-equilibrium work associated with all possibilities with dynamics  $x_t$  starting at  $x_0$  and ending at  $x_T$  (for an example, see ref 53). Therefore, eq 46 is the detailed Jarzynski relation. Noticing that

$$p(x_T, T|x_0, 0) = \int_{C\{x_0, 0; x_T, T\}} dp[x_t] \quad (49)$$

From eq 46, we have

$$\frac{p(x_0, T|x_T, 0)}{p(x_T, T|x_0, 0)} = \frac{\mu(x_0)}{\mu(x_T)} \mathbb{E}_{x_0 \rightarrow x_T} [e^{-\beta W}] \quad (50)$$

## APPENDIX B. COMPUTATION OF $\hat{T}$ AND $\hat{M}$ DIRECTLY FROM TRAJECTORIES

To show how to compute  $\hat{T}$  and  $\hat{M}$  directly from trajectories, let us start by denoting the first hitting time of set  $A$  starting at time  $t$  by  $h_t(A)$ . In addition, we define  $\hat{q}_j^-(i) = P(X_t = i | \hat{X}_t = j)$  and always assume  $t = mT$ . Then, due to the Bayes' theorem, we have

$$\hat{q}_j^-(i) = \frac{P(\hat{X}_t = j | X_t = i) P(X_t = i)}{P(\hat{X}_t = j)} = \frac{q_j^-(i) \mu(i)}{\hat{\mu}_j} \quad (51)$$

Then

$$\begin{aligned} & P[h_t(C_k) < h_t(\cup_{l \neq k} C_l) | \hat{X}_t = j] \\ &= \sum_{i=1}^N \frac{P[h_t(C_k) < h_t(\cup_{l \neq k} C_l), X_t = i, \hat{X}_t = j]}{P(\hat{X}_t = j)} \\ &= \sum_{i=1}^N P[h_t(C_k) < h_t(\cup_{l \neq k} C_l) | X_t = i, \hat{X}_t = j] \hat{q}_j^-(i) \\ &= \sum_{i=1}^N P[h_t(C_k) < h_t(\cup_{l \neq k} C_l) | X_t = i] \hat{q}_j^-(i) \\ &= \sum_{i=1}^N q_k^+(i) \frac{q_j^-(i) \mu(i)}{\hat{\mu}_j} = \frac{\langle q_j^-, q_k^+ \rangle_\mu}{\hat{\mu}_j} = \hat{M}_{jk} \end{aligned} \quad (52)$$

where the third equation holds because of the Markovianity of process  $X_t$  and the fourth equation is due to the definition of the forward committor. The interpretation of the result is that when starting with the milestone process in state  $j$ , we have to determine the fraction of all trajectories that hit core set  $C_k$  first from all core sets to estimate  $\hat{M}_{jk}$ .

To see the same for  $\hat{T}$ , we first need

$$\begin{aligned} & P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l) | X_t = i] \\ &= \sum_{l=1}^N \frac{P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l), X_{t+T} = l, X_t = i]}{P(X_t = i)} \\ &= \sum_{l=1}^N P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l) | X_{t+T} = l, X_t = i] \\ &\quad \times P(X_{t+T} = l | X_t = i) \\ &= \sum_{l=1}^N P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l) | X_{t+T} = l] \\ &\quad \times P(X_{t+T} = l | X_t = i) = \sum_{l=1}^N q_k^+(l) P_{il} \end{aligned}$$

We used the Markovianity of  $X_T$  in the third equation and the time-homogeneity in the fourth equation. Therefore, following the same procedure as before, we get

$$\begin{aligned} & P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l) | \hat{X}_t = j] \\ &= \sum_{i=1}^N \sum_{l=1}^N q_k^+(l) P_{il} \frac{q_j^-(i) \mu(i)}{\hat{\mu}_j} = \frac{\langle q_j^-, P q_k^+ \rangle_\mu}{\hat{\mu}_j} = \hat{T}_{jk} \end{aligned} \quad (53)$$

Thus, when starting with the milestone process in state  $j$  at some time  $t$ , we have to determine the fraction of all trajectories of length at least  $T$  that hit core set  $C_k$  first of all of the core sets to estimate  $\hat{T}_{jk}$ .

## APPENDIX C. PROVING OF THE IDENTITY $\hat{P} = \hat{T} \hat{M}^{-1}$

In this section, we prove the identity  $\hat{P} = \hat{T} \hat{M}^{-1}$ . Starting from eq 53,

$$\begin{aligned} \hat{T}_{jk} &= P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l) | \hat{X}_t = j] \\ &= \frac{P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l), \hat{X}_t = j]}{P(\hat{X}_t = j)} \\ &= \sum_l \frac{P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l), \hat{X}_{t+T} = l, \hat{X}_t = j]}{P(\hat{X}_t = j)} \\ &= \sum_l P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l) | \hat{X}_{t+T} = l, \hat{X}_t = j] \hat{P}_{jl} \\ &= \sum_l P[h_{t+T}(C_k) < h_{t+T}(\cup_{l \neq k} C_l) | \hat{X}_{t+T} = l] \hat{P}_{jl} \\ &= \sum_l \hat{P}_{jl} \hat{M}_{lk} \end{aligned}$$

This proves the identity  $\hat{P} = \hat{T} \hat{M}^{-1}$ . Notice that in the fifth equation, we assume the Markovianity of the milestone process  $\hat{X}_m$ . In the sixth equation, we use the result from eq 52 and the time-homogeneity of the milestone process.

## AUTHOR INFORMATION

### Corresponding Authors

\*E-mail: wang\_han@iapcm.ac.cn.

\*E-mail: schuette@zib.de.

### Notes

The authors declare no competing financial interest.

## ■ ACKNOWLEDGMENTS

The computation of this work was performed at the North-German Supercomputing Alliance (HLRN). The authors are grateful for this support. H.W. acknowledges support from the National High Technology Research and Development Program of China under Grant 2015AA011201.

## ■ REFERENCES

- (1) Bohr, H.; Bohr, J. *Phys. Rev. E* **2000**, *61*, 4310.
- (2) Bohr, H.; Bohr, J. *Bioelectromagnetics* **2000**, *21*, 68–72.
- (3) de Pomerai, D.; Daniells, C.; David, H.; Allan, J.; Duce, I.; Mutwakil, M.; Thomas, D.; Sewell, P.; Tattersall, J.; Jones, D.; Candido, P. *Nature* **2000**, *405*, 417–418.
- (4) de Pomerai, D. I.; Smith, B.; Dawe, A.; North, K.; Smith, T.; Archer, D. B.; Duce, I. R.; Jones, D.; Candido, E. P. M. *FEBS Lett.* **2003**, *543*, 93–97.
- (5) Mancinelli, F.; Caraglia, M.; Abbruzzese, A.; d'Ambrosio, G.; Massa, R.; Bismuto, E. *J. Cell. Biochem.* **2004**, *93*, 188–196.
- (6) Inskip, P. D.; Tarone, R. E.; Hatch, E. E.; Wilcosky, T. C.; Shapiro, W. R.; Selker, R. G.; Fine, H. A.; Black, P. M.; Loeffler, J. S.; Linet, M. S. *N. Engl. J. Med.* **2001**, *344*, 79–86.
- (7) Bekard, I.; Dunstan, D. E. *Soft Matter* **2014**, *10*, 431–437.
- (8) Budi, A.; Legge, F.; Treutlein, H.; Yarovsky, I. *J. Phys. Chem. B* **2005**, *109*, 22641–22648.
- (9) Budi, A.; Legge, F.; Treutlein, H.; Yarovsky, I. *J. Phys. Chem. B* **2007**, *111*, 5748–5756.
- (10) Budi, A.; Legge, F. S.; Treutlein, H.; Yarovsky, I. *J. Phys. Chem. B* **2008**, *112*, 7916–7924.
- (11) Astrakas, L. G.; Gousias, C.; Tzaphlidou, M. J. *J. Appl. Phys.* **2012**, *111*, 074702–074702.
- (12) Damm, M.; Nusshold, C.; Cantillo, D.; Rechberger, G. N.; Gruber, K.; Sattler, W.; Kappe, C. O. *J. Proteomics* **2012**, *75*, 5533–5543.
- (13) English, N.; Solomentsev, G.; O'Brien, P. J. *Chem. Phys.* **2009**, *131*, 035106.
- (14) Solomentsev, G.; English, N.; Mooney, D. J. *Comput. Chem.* **2012**, *33*, 917–923.
- (15) Wang, H.; Schütte, C.; Ciccotti, G.; Delle Site, L. *J. Chem. Theory Comput.* **2014**, *10*, 1376–1386.
- (16) Schütte, C.; Sarich, M. *Metastability and Markov State Models in Molecular Dynamics: Modeling, Analysis, Algorithmic Approaches*; American Mathematical Society: Providence, RI, 2013; Vol. 24, Courant lecture notes.
- (17) Prinz, J.; Wu, H.; Sarich, M.; Keller, B.; Senne, M.; Held, M.; Chodera, J.; Schütte, C.; Noé, F. *J. Chem. Phys.* **2011**, *134*, 174105.
- (18) Bowman, G. R.; Pande, V. S.; Noé, F., Eds. *Advances in Experimental Medicine and Biology: An Introduction to Markov State Models and Their Application to Long Timescale Molecular Simulation*; Springer: Dordrecht, The Netherlands, 2014; Vol. 797.
- (19) Senne, M.; Trendelkamp-Schroer, B.; Mey, A. S. J. S.; Schütte, C.; Noé, F. *J. Chem. Theory Comput.* **2012**, *8*, 2223–2238.
- (20) Beauchamp, K. A.; Bowman, G. R.; Lane, T. J.; Maibaum, L.; Haque, I. S.; Pande, V. S. *J. Chem. Theory Comput.* **2011**, *7*, 3412–3419.
- (21) Schütte, C.; Noé, F.; Lu, J.; Sarich, M.; Vanden-Eijnden, E. *J. Chem. Phys.* **2011**, *134*, 204105.
- (22) Sarich, M.; Noé, F.; Schütte, C. *Multiscale Model. Simul.* **2010**, *8*, 1154–1177.
- (23) Djurdjevac, N.; Sarich, M.; Schütte, C. *Multiscale Model. Simul.* **2012**, *10*, 61–81.
- (24) Noé, F.; Schütte, C.; Vanden-Eijnden, E.; Reich, L.; Weikl, T. R. *Proc. Natl. Acad. Sci. U.S.A.* **2009**, *106*, 19011–19016.
- (25) Kohlhoff, K. J.; Shukla, D.; Lawrenz, M.; Bowman, G. R.; Konerding, D. E.; Belov, D.; Altman, R. B.; Pande, V. S. *Nat. Chem.* **2014**, *6*, 15–21.
- (26) Keller, B. G.; Prinz, J.-H.; Noé, F. *Chem. Phys.* **2012**, *396*, 92–107.
- (27) Prinz, J.-H.; Keller, B. G.; Noé, F. *Phys. Chem. Chem. Phys.* **2011**, *13*, 16912–16927.
- (28) Pande, V.; Beauchamp, K.; Bowman, G. *Methods* **2010**, *52*, 99–105.
- (29) Sarich, M.; Banisch, R.; Hartmann, C.; Schütte, C. *Entropy* **2013**, *16*, 258–286.
- (30) Latorre, J. C.; Metzner, P.; Hartmann, C.; Schütte, C. *Communications in Mathematical Sciences* **2011**, *9*, 1051–1072.
- (31) Floquet, G. *Annales Scientifiques de l'École Normale Supérieure* **1883**, *12*, 47–82.
- (32) Noé, F.; Wu, H.; Prinz, J.-H.; Plattner, N. *J. Chem. Phys.* **2013**, *139*, 184114.
- (33) Buchete, N. V.; Hummer, G. *J. Phys. Chem. B* **2008**, *112*, 6057–6069.
- (34) Sarich, M.; Schütte, C. *J. Comput. Dyn.* 2014, submitted.
- (35) Metzner, P.; Schütte, C.; Vanden-Eijnden, E. *Multiscale Model. Simul.* **2009**, *7*, 1192–1219.
- (36) Djurdjevac, N.; Sarich, M.; Schütte, C. *On Markov State Models for Metastable Processes*; Proceedings of the International Congress of Mathematicians; New Delhi, India, **2010**; pp 3105–3131.
- (37) Chodera, J.; Singhal, N.; Pande, V.; Dill, K.; Swope, W. *J. Chem. Phys.* **2007**, *126*, 155101.
- (38) Foloppe, N.; MacKerell, A. D. *J. Comput. Chem.* **2000**, *21*, 86–104.
- (39) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (40) MacKerell, A. D.; Feig, M.; Brooks, C. L. *J. Comput. Chem.* **2004**, *25*, 1400–1415.
- (41) Pronk, S.; Páll, S.; Schulz, R.; Larsson, P.; Bjelkmar, P.; Apostolov, R.; Shirts, M.; Smith, J.; Kasson, P.; van der Spoel, D.; Hess, B.; Lindahl, E. *Bioinformatics* **2013**, *1*–10.
- (42) Bjelkmar, P.; Larsson, P.; Cuendet, M. A.; Hess, B.; Lindahl, E. *J. Chem. Theory Comput.* **2010**, *6*, 459–466.
- (43) Parrinello, M.; Rahman, A. *J. Appl. Phys.* **1981**, *52*, 7182.
- (44) Darden, T.; York, D.; Pedersen, L. *J. Chem. Phys.* **1993**, *98*, 10089.
- (45) Essmann, U.; Perera, L.; Berkowitz, M.; Darden, T.; Lee, H.; Pedersen, L. *J. Chem. Phys.* **1995**, *103*, 8577.
- (46) Wang, H.; Dommert, F.; Holm, C. *J. Chem. Phys.* **2010**, *133*, 034117.
- (47) Hess, B.; Bekker, H.; Berendsen, H.; Fraaije, J. *J. Comput. Chem.* **1997**, *18*, 1463–1472.
- (48) Miyamoto, S.; Kollman, P. *J. Comput. Chem.* **2004**, *13*, 952–962.
- (49) Deuffhard, P.; Weber, M. *Linear Algebra and its Applications*, Special Issue on Matrices and Mathematical Biology; **2005**; Vol. 161, 398.
- (50) Prinz, J.-H.; Held, M.; Smith, J. C.; Noé, F. *Multiscale Model. Simul.* **2011**, *9*, 545.
- (51) Wang, H.; Hartmann, C.; Schütte, C. *Mol. Phys.* **2013**, *111*, 3555–3564.
- (52) Schütte, C.; Nielsen, A.; Weber, M. *Mol. Phys.* **2015**, *113*, 69–78.
- (53) Seifert, U. *Rep. Prog. Phys.* **2012**, *75*, 126001.