

Article

Markov State Models for Rare Events in Molecular Dynamics

Marco Sarich^{1,*}, Ralf Banisch¹, Carsten Hartmann¹ and Christof Schütte^{1,2}

¹ Department of Mathematics and Computer Science, Freie Universität Berlin, Arnimallee 6, 14195 Berlin, Germany

² Zuse Institute Berlin, Takustr. 7, 14195 Berlin, Germany

* Author to whom correspondence should be addressed; sarich@math.fu-berlin.de

Version June 2, 2013 submitted to *Entropy*. Typeset by *LaTeX* using class file *mdpi.cls*

Abstract: Rare but important transition events between long lived states are a key feature of many molecular systems. In many cases the computation of rare event statistics by direct molecular dynamics (MD) simulations is infeasible even on the most powerful computers because of the immensely long simulation timescales needed. Recently a technique for spatial discretization of the molecular state space designed to help overcome such problems, so-called Markov State Models (MSMs), has attracted a lot of attention. We review the theoretical background and algorithmic realization of MSMs and illustrate their use by some numerical examples. Furthermore we introduce a novel approach to using MSMs for the efficient solution of optimal control problems that appear in applications where one desires to optimize molecular properties by means of external controls.

Keywords: rare events; Markov state models; long timescales; optimal control

1. Introduction

Stochastic processes are widely used to model physical, chemical or biological systems. The goal is to approximately compute interesting properties of the system by analyzing the stochastic model. There are mainly two options for performing this analysis: (1) Direct or accelerated sampling of the process and (2) the construction of a discrete coarse grained model of the system. In a sampling approach, one tries to generate a statistically significant amount of events that characterize the property of the system one is interested in. For this purpose, computer simulations of the model are a powerful tool. For example, an event could refer to the transition between two well-defined macroscopic states of the system. In chemical applications such transitions can often be interpreted as reactions, or in the context

21 of a molecular system as conformational changes. Interesting properties are e.g. average waiting times
22 for such reactions or conformational changes and along which pathways the transitions typically occur.
23 The problem with a direct sampling approach is that many interesting events are so called rare events.
24 Therefore the computational effort for generating sufficient statistics for reliable estimates is very high,
25 and, particularly if the state space is continuous and high dimensional, estimation by direct numerical
26 simulation is infeasible. Accelerated sampling tries to overcome this problem but cannot be applied in
27 general.

28 Available techniques for rare event simulations in continuous state space are discussed in [1]. In this
29 article, we will discuss approach (2) to the estimation of rare event statistics via *discretization* of the
30 state space of the system under consideration. That is, instead of dealing with the computation of rare
31 events for the original, continuous process, we will approximate them by a so-called Markov State Model
32 (MSM) with *discrete* finite state space. The reason is that for such a discrete model one can numerically
33 compute many interesting properties without simulation, mostly by solving linear systems of equations
34 as in discrete transition path theory (TPT) [2]. We will see that this approach, called Markov State
35 Modelling, *avoids* the combinatorial explosion of the number of discretization elements with increasing
36 size of the molecular system in contrast to other methods for spatial discretization.

37 The actual construction of an MSM requires to sample certain transition probabilities of the
38 underlying dynamics between sets. The idea is (1) to choose the sets such that the sampling effort is much
39 lower than the direct estimation of the rare events under consideration, and (2) to compute all interesting
40 quantities for the MSM from its transition matrix, cf. [2,3]. There are many examples for the successful
41 application of this strategy. In [4], for example, it was used to compute dominant folding pathways for
42 the PinWW domain in explicit solvent. However, we have to make sure that the Markov State Model
43 approximates the original dynamics well enough. For example, the MSM should correctly reproduce
44 the timescales of the processes of interest. These approximation issues have been discussed since more
45 than a decade now [5,6]; in this article we will review the present state of research on this topic. In
46 the algorithmic realization of Markov State Modelling for realistic molecular systems the transition
47 probabilities and the respective statistical uncertainties are estimated from short MD trajectories only, cf.
48 [7]. This makes Markov State Modelling applicable to many different molecular systems and processes,
49 cf. [8–13].

50 In the first part of this article we will discuss the approximation quality of two different types of
51 Markov State Models that are defined with respect to a full partition of state space or with respect to
52 so-called core sets. We will also discuss the algorithmic realization of MSMs and provide references
53 to the manifold of realistic applications to molecular systems in equilibrium that are available in the
54 literature today.

55 The second part will show how to use MSMs for optimizing particular molecular properties. In this
56 type of application one wants to steer the molecular system at hand by external controls in a way such
57 that a pre-selected molecular property is optimized (minimized or maximized). That is, one wants to
58 compute a specific external control from a family of admissible controls that optimizes the property of
59 interest under certain side conditions. The property to be optimized can be quite diverse: For example,
60 it can be (1) the population of a certain conformation that one wants to maximize under a side condition
61 that limits the total work done by the external control or (2) the mean first passage time to a certain

62 conformation that one wants to minimize (in order to speed up a rare event) but under the condition
 63 that one can still safely estimate the mean first passage time of the uncontrolled system. The theoretical
 64 background of case (1) has been considered in [14], for example, and of case (2) in [1,15]. There one
 65 finds the mathematical problem that has to be solved in order to compute the optimal control. Here we
 66 will demonstrate that one can use MSMs for the efficient solution of such a mathematical problem (for
 67 both cases). We will see that the spatial discretization underlying an MSM turns the high-dimensional
 68 continuous optimal control problem into a rather low-dimensional discrete optimal control problem of
 69 the same form that can be solved efficiently. Based on these insights, MSM discretization yields an
 70 efficient algorithm for solving the optimal control problem whose performance we will outline in some
 71 numerical examples including an application to Alanine dipeptide.

72 2. MSM Construction

73 Let $(X_t)_{t \geq 0}$ be a time-continuous Markov process on a continuous state space E , e.g. $E \subset \mathbb{R}^d$.
 74 That is, X_t is the state of the molecular system at time t as resulting from any usually used from
 75 of molecular dynamics simulation, be it based on Newtonian dynamics with thermostats or resulting
 76 from Langevin dynamics or other diffusion molecular dynamics models. The idea of Markov State
 77 Modelling is to derive a Markov chain $(\hat{X}_k)_{k \in \mathbb{N}}$ on a finite and preferably small state space $\hat{E} = \{1, \dots, n\}$
 78 that models characteristic dynamics of the continuous process (X_t) . For example, in molecular
 79 dynamics applications such characteristic dynamics could refer to protein folding processes [16,17],
 80 conformational rearrangements between native protein substates [18,19], or ligand binding processes
 81 [20]. Since the approximating Markov chain $(\hat{X}_k)_{k \in \mathbb{N}}$ lives on a finite state space, the construction of an
 82 MSM boils down to the computation of its transition matrix P

$$P_{ij} = \mathbb{P}[\hat{X}_{k+1} = j | \hat{X}_k = i]. \quad (1)$$

83 The main benefit is that for a finite Markov chain one can compute many interesting dynamical
 84 properties directly from its transition matrix, e.g. timescales and metastability in the system [5,21,22], a
 85 hierarchy of important transition pathways [2], or mean first passage times between selected states. With
 86 respect to na MSM, these computations should be used afterwards to answer related questions for the
 87 original continuous process. To do this we must be able to link the states of the Markov chain back to
 88 spatial information of the original process and the approximation of the process (X_t) by the MSM must
 89 be valid in some sense.

Having this in mind the first natural idea is to let the states of an MSM correspond to sets $A_1, \dots, A_n \subset E$ in continuous state space that form a full partition, i.e.

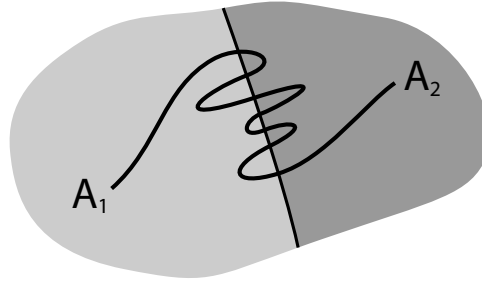
$$A_i \cap A_j = \emptyset \text{ for } i \neq j, \quad \bigcup_{i=1}^n A_i = E. \quad (2)$$

90 Typical choices for such sets are box discretizations or voronoi tessellations [23]. For such a full partition
 91 it is trivial to also define a corresponding discretized process by the original switching dynamics between
 92 the sets. For a given lag time $\tau > 0$, we can define the index process

$$\tilde{X}_k = i \Leftrightarrow X_{k\tau} \in A_i. \quad (3)$$

93 It is well known that this process is not Markovian, mainly due to the so called recrossing problem.
 94 It refers to the fact that the original process typically crosses the boundary between two sets A_i and A_j
 95 several times when transitions take place, as illustrated in Fig. 1. This results in cumulative transitions
 96 between indices i and j for the index process, that is, a not memoryless transition behavior.

Figure 1. Cumulative transitions between two sets along boundaries are typical.



97 The non-Markovianity of the index process is often seen as a problem in Markov State Modeling
 98 because many arguments assume that \tilde{X}_k is a Markov process. In this article, we will *not* make this
 99 assumption. We interpret the process (\tilde{X}_k) as a tool to construct the following transition matrix P^τ

$$P_{ij}^\tau = \mathbb{P}[\tilde{X}_{k+1} = j | \tilde{X}_k = i] = \mathbb{P}[X_{(k+1)\tau} \in A_j | X_{k\tau} \in A_i] \quad (4)$$

100 and hence the MSM as the Markov chain $(\hat{X}_k)_{k \in \mathbb{N}}$ associated with this transition matrix. From above
 101 it is clear that in general we have $\hat{X}_k \neq \tilde{X}_k$ and in [24] it was analyzed how these two processes
 102 relate in terms of density propagation. In the following, we will show under which assumptions and
 103 in which sense the MSM (\hat{X}_k) will be a good approximation of the original dynamics given by (X_t) . For
 104 convenience we will usually write $P^\tau \equiv P$ and leave the τ -dependence implicit.

105 3. Analytical Results

In order to compare the MSM to the continuous process we introduce one of the key objects for our analysis, the *transfer operator* of a Markov process. We assume that the Markov process (X_t) has a unique, positive invariant probability measure μ and that it is time-reversible. Then, for any time-step $t \geq 0$ we define the transfer operator T_t via the property

$$\int_A T_t v(y) \mu(dy) = \int_E v(x) p(t, x, A) \mu(dx) \quad \text{for all measurable } A \quad (5)$$

as an operator $T_t : L^2(\mu) \rightarrow L^2(\mu)$. Here, $p(t, x, A) = \mathbb{P}[X_t \in A | X_0 = x]$ defines the transition probability measure and $L^2(\mu)$ denotes the Hilbert space of functions v with

$$\int_E v(y)^2 \mu(dy) \leq \infty \quad (6)$$

and the scalar product

$$\langle v, w \rangle = \int_E v(y) w(x) \mu(dy). \quad (7)$$

Note that T_t is nothing else than the propagator of densities under the dynamics, but the densities are understood as densities with respect to the measure μ . That is, if the Markov process is initially distributed according to

$$\mathbb{P}[X_0 \in A] = \int_A v_0(x) \mu(dx), \quad (8)$$

its probability distribution at time t is given by

$$\mathbb{P}[X_t \in B] = \int_B v_t(x) \mu(dx), \quad v_t = T_t v_0. \quad (9)$$

The benefit of working with μ -weighted densities is that the transfer operator T_t becomes essentially self-adjoint on $L^2(\mu)$ for all cases of molecular dynamics satisfying some form of detailed balance condition. Hence, it has real eigenvalues and orthogonal eigenvectors with respect to (7) (or at least the dominant spectral elements are real-valued). Moreover, the construction of an MSM can be seen as a projection of the transfer operator [25]. Assume Q is an orthogonal projection in $L^2(\mu)$ onto an n -dimensional subspace $D \subset L^2(\mu)$ with $\mathbb{1} \in D$, and χ_1, \dots, χ_n is a basis of D . Then, the so called projected transfer operator $QT_\tau Q : D \rightarrow D$ has the matrix representation

$$P_Q = PM^{-1}, \quad (10)$$

with the non-negative, invertible mass matrix $M \in \mathbb{R}^{n,n}$ with entries

$$M_{ij} = \frac{\langle \chi_i, \chi_j \rangle}{\langle \chi_i, \mathbb{1} \rangle}. \quad (11)$$

The matrix $P \in \mathbb{R}^{n,n}$ is also non-negative and has entries

$$P_{ij} = \frac{\langle \chi_i, T_\tau \chi_j \rangle}{\langle \chi_i, \mathbb{1} \rangle}. \quad (12)$$

Full Partition MSM. If we choose $\chi_i = \mathbb{1}_{A_i}$ to be the characteristic function of set A_i for $i = 1, \dots, n$, one can easily check that we get $M = I$ to be the identity matrix and

$$P_{ij} = \mathbb{P}_\mu[X_\tau \in A_j | X_0 \in A_i] \quad (13)$$

106 as in (4). The subscript μ shall indicate that $X_0 \sim \mu$. So the transition probabilities are evaluated along
107 equilibrium paths.

The previously constructed transition matrix of the MSM based on a full partition can be interpreted as a projection onto a space of densities which are constant on the partitioning sets. This interpretation of an MSM is useful since it allows to analyze its approximation quality. For example, in [25,26] it is proven that we can reproduce an eigenvalue λ of a self-adjoint transfer operator T_t by the MSM by choosing the subspace appropriately. That is, if u is a corresponding normalized eigenvector, Q the orthogonal projection to a subspace D with $\mathbb{1} \in D$, then there exists an eigenvalue $\hat{\lambda}$ of the projected transfer operator QT_tQ with

$$|\lambda - \hat{\lambda}| \leq \lambda_1 \delta (1 - \delta^2)^{-\frac{1}{2}},$$

108 where $\lambda_1 < 1$ is the largest non-trivial eigenvalue of T_t and $\delta = \|u - Qu\|$.

In particular, for $\delta \leq \frac{3}{4}$ one can simplify the equation to

$$|\lambda - \hat{\lambda}| \leq 2\lambda_1\delta. \quad (14)$$

109

An eigenvalue λ_i of the transfer operator directly relates to an implied timescales \mathcal{T}_i of the system via

$$\mathcal{T}_i = -\frac{\tau}{\log(\lambda_i)}. \quad (15)$$

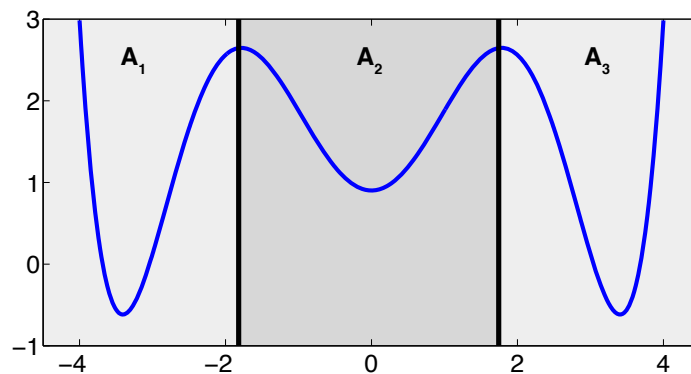
110 So the transition matrix (4) that we construct from transitions between the sets A_1, \dots, A_n will generate
 111 a Markov chain that will reproduce the original timescales well if the partitioning sets are chosen such
 112 that the corresponding eigenvectors are almost constant on these sets. In this case $\delta = \|u - Qu\|$, that is
 113 the approximation error of the eigenvector by a piecewise constant function on the sets will be small.

The projection error δ depends on our choice of the discretizing sets. As an example let us consider a diffusion in the potential that is illustrated in Fig. 2, that is, the reversible Markov process given by the stochastic differential equation

$$dX_t = -\nabla V(X_t)dt + \sqrt{2\varepsilon}dB_t, \quad (16)$$

114 where V is the potential, B_t denotes a Brownian motion and $\varepsilon > 0$.

Figure 2. A potential with three wells and a choice of 3 sets A_1, A_2, A_3 .



The figure also shows a choice of three sets that form a full partition of state space. The computation of the transition matrix (4) for $\sigma = 0.7$ and a lag time $\tau = 1$ yields

$$P_Q = P = \begin{pmatrix} 0.9877 & 0.0123 & 0.0000 \\ 0.0420 & 0.9160 & 0.0419 \\ 0.0000 & 0.0123 & 0.9877 \end{pmatrix}$$

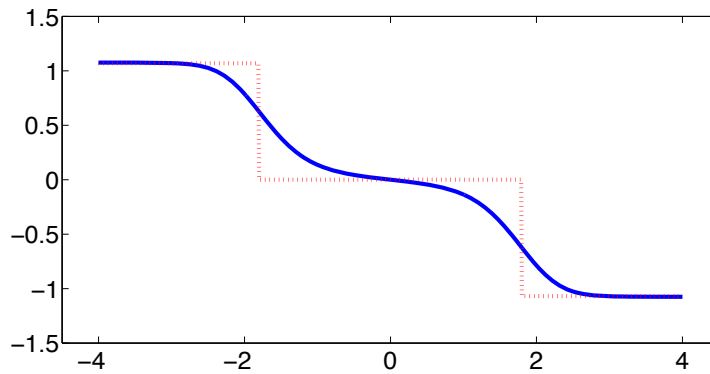
115 that has three eigenvalues $\lambda_0 = 1, \lambda_1 = 0.9877, \lambda_2 = 0.9037$. The following table shows the two
 116 resulting implied timescales (15) in comparison to the timescales of the original system.

	T_1	T_2
original	103.7608	11.9566
full partition 3 sets	80.6548	9.8784

117

118 As one can see, the timescales are strongly underestimated. This is a typical phenomenon. From a
 119 statistical point of view, the recrossing problem will lead to cumulatively appearing transition counts
 120 when one computes the transition probabilities $\mathbb{P}_\mu[X_\tau \in A_j | X_0 \in A_i]$ from a trajectory (X_t) , as
 121 discussed above. Therefore on average transitions between sets seem to become too likely and hence
 122 the processes in the coarse grained system get accelerated. We have seen in (14) that this cannot happen
 123 if the associated eigenvectors can be approximated well by the subspace that corresponds to the MSM.
 124 Fig. 3 shows the first non-trivial eigenvector u_1 belonging to the timescale $T_1 = 103.7608$ and its
 125 best-approximation by a step function.

Figure 3. The first non-trivial eigenvector u_1 (solid blue) and its projection Qu_1 (dashed red) onto step functions that are constant on A_1, A_2, A_3 .



126 The eigenvector is indeed almost constant in the vicinity of the wells, but within the transition region
 127 between the wells the eigenvector is varying and the approximation by a step function is not accurate.
 128 So we have two explanations why the main error is introduced in the region close to shared boundaries
 129 of neighboring sets: (1) because of recrossing issues and (2) because of the main projection error of the
 130 associated eigenvector. Of course, one solution would be an adaptive refinement of the discretization,
 131 that is, one could choose a larger number of smaller sets such that the eigenvector is better approximated
 132 by a step function on these sets. In the following section, we will present an alternative solution for
 133 overcoming the recrossing problem and reducing the projection error without refining the discretization.

134 4. The Core Set Approach

From (10) we know how to compute a matrix representation for a projected transfer operator for an arbitrary subspace $D \subset L^2(\mu)$. For a given basis χ_1, \dots, χ_n we have to compute (11) and (12), so

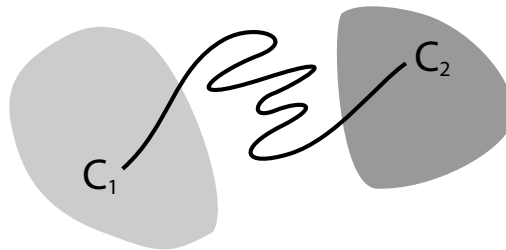
$$M_{ij} = \frac{\langle \chi_i, \chi_j \rangle}{\langle \chi_i, \mathbf{1} \rangle}, \quad P_{ij} = \frac{\langle \chi_i, T_\tau \chi_j \rangle}{\langle \chi_i, \mathbf{1} \rangle}. \quad (17)$$

135 In general, the evaluation of these scalar products for arbitrary basis functions is a non-trivial task. On
 136 the other hand, we have seen that for characteristic functions $\chi_i = \mathbf{1}_{A_i}$ on a full partition we do not
 137 have to compute the scalar products numerically since the matrix entries have a stochastic interpretation
 138 in terms of transition probabilities between sets (13). This means they can be directly estimated from

139 a trajectory of the process which is a strong computational advantage, particularly in high dimensional
140 state spaces.

141 Now, the question is if there is another basis than characteristic functions that a) is more adapted
142 to the eigenvectors of the transfer operator, and b) still leads to a probabilistic interpretation of the
143 matrix entries (17) such that scalar products never have to be computed. The basic idea is to stick to a
144 set-oriented definition of the basis, but to relax the full partition constraint. We will define our basis with
145 respect to so called *core sets* $C_1, \dots, C_n \subset E$ that are still disjoint, so $C_i \cap C_j = \emptyset$, but they do not have to
146 form a full partition. Figure 4 suggests that this could lead to a reduction of the recrossing phenomenon
147 since the sets do not share boundaries anymore.

Figure 4. Core sets do not have to share boundaries anymore. This can reduce the recrossing effect.



148 Now, we use the core sets to define our basis functions χ_1, \dots, χ_n . Assume T_τ is again a self-adjoint
149 transfer operator and consider n core sets C_1, \dots, C_n . For every i , take the committor function χ_i of the
150 process with respect to core set C_i , that is, $\chi_i(x)$ denotes the probability to hit the core set C_i next rather
151 than the other core sets when starting the process in x . If we now study the the projection Q onto the
152 space spanned by these committor functions, the two following properties hold [25,27].

(P1) The matrices M and P in (10) can be written as

$$M_{ij} = \mathbb{P}_\mu[\tilde{X}_k^+ = j | \tilde{X}_k^- = i], \quad P_{ij} = \mathbb{P}_\mu[\tilde{X}_{k+1}^+ = j | \tilde{X}_k^- = i], \quad (18)$$

153 where (\tilde{X}_k^+) and (\tilde{X}_k^-) are forward and backward milestone processes [25,28], that is, $\tilde{X}_k^- = i$
154 if the process came at time $t = k\tau$ last from core set C_i and $\tilde{X}_k^+ = j$ if the process went next to
155 core set C_j after time $t = k\tau$.

156 (P2) Let u_i be an eigenvector of T_τ that is almost constant on the core sets. Let the region $C = E \setminus \bigcup_i C_i$
157 that is not assigned to a core set be left quickly enough, so $\mathbb{E}_x[\tau(C^c)] \ll \mathcal{T}_i$ for all $x \in C$, where \mathcal{T}_i
158 is the timescale associated with u_i and $\mathbb{E}_x[\tau(C^c)]$ is the expected hitting time of $C^c = \bigcup_i C_i$ when
159 starting in $x \in C$. Then, $\|u_i - Qu_i\|$ is small, so the committor approximation to the eigenvector
160 is accurate.

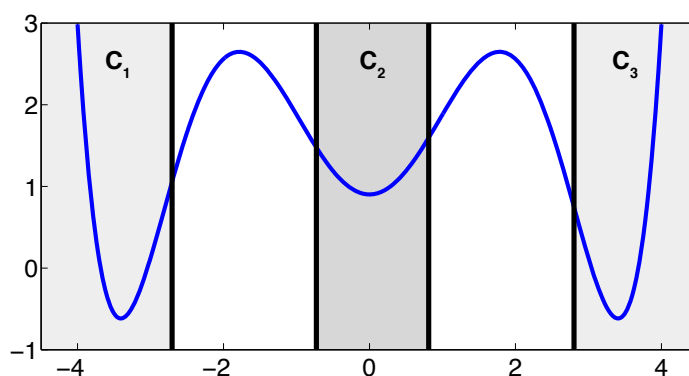
161 The message behind (P1) is that it is possible to relax the full partition constraint and use a core set
162 discretization that does not cover the whole state space. We can still define a basis for a projection of
163 the transfer operator that leads to a matrix representation that can be interpreted in terms of transition
164 probabilities.

165 **Important remark:** The construction of the projection onto the committors is only necessary for
 166 theoretical purposes. In practice, neither the committor functions, nor scalar products between the
 167 committors have to be computed numerically, since the matrix entries of M and P can be estimated
 168 from trajectories again.

169

170 Property (P2) yields that the relaxation of the full partition constraint should also lead to an
 171 improvement of the MSM if the region C between the core sets is typically left on a faster timescale
 172 than the processes of interest take place. Let us get back to the example from above. We will see that we
 173 can achieve a strong improvement of the approximation by simply excluding a small part of state space
 174 from our discretization. In Figure 5 we have turned our initial full partition into a core set discretization
 175 by removing parts of the transition region between the wells.

Figure 5. Excluding a small region of state space from the sets A_1, A_2, A_3 as in Fig. 2 to form core sets C_1, C_2, C_3 that do not share boundaries anymore.



176 The matrix $P_Q = PM^{-1}$ that represents the projection $QT_\tau Q$ of the transfer operator onto the
 177 committor space associated with the core sets is given by

$$P_Q = \begin{pmatrix} 0.9897 & 0.0103 & 0.0000 \\ 0.0352 & 0.9298 & 0.0351 \\ 0.0000 & 0.0103 & 0.9897 \end{pmatrix}$$

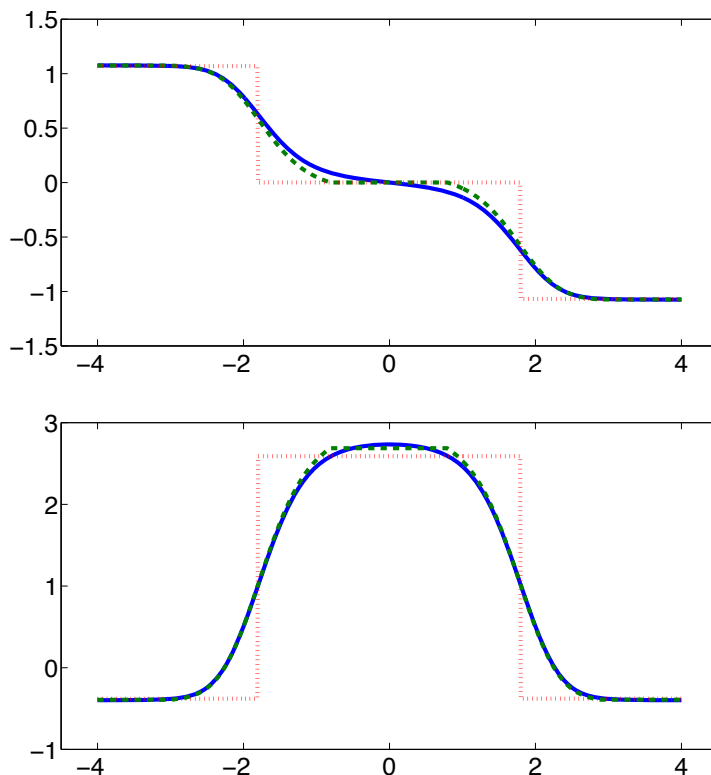
178 Comparing to the MSM for the full partition one can see that transitions between indices i and j , $i \neq j$ are
 179 less likely. As the following table shows this leads to a far more accurate reproduction of the timescales
 180 in the system.

	T_1	T_2
original	103.7608	11.9566
3 core sets	100.8066	11.9145
full partition 3 sets	80.6548	9.8784

182 From the discussion above this has to be expected because the eigenvectors are almost constant in
 183 the vicinity of the wells and we removed a part of state space from the discretization that is typically

184 left quickly compared to the timescales T_1 and T_2 . So, the committor functions should deliver a good
 185 approximation of the first two eigenvectors. Figure 6 underlines this theoretical result.

Figure 6. Upper panel: The first non-trivial eigenvector u_1 (solid blue) and its projection $Q_f u_1$ (finely dashed red) onto stepfunctions (full partition) and its projection $Q_c u_1$ (dashed green) onto committors (core sets). Lower panel: The same plot for the second non-trivial eigenvector u_2 .



186 5. Practical Considerations and MD Applications

In the previous sections we have interpreted the construction of an MSM as projection of the dynamics onto some finite dimensional ansatz space. We have discussed two types of spaces that both have been defined on the basis of a set discretization. First, we chose a full partition of state space and the associated space of step functions, and second we analyzed a discretization by core sets and the associated space spanned by committor functions. These two methods have the advantage that the resulting projections lead to transition matrices for the MSM with entries that are given in terms of transition probabilities between the sets. That is, one can compute estimates for the transition matrices from simulation data. This is an important property for practical applications because it means that we never need to compute committor functions, or scalar products between committors or step functions. We rather generate trajectories x_0, x_1, \dots, x_N of the process (X_t) , let us say for a time step $h > 0$, so

$x_i = X_{hi}$. For example, we can then define for a full partition A_1, \dots, A_m and a lag time $\tau = nh$ the discrete trajectory $s_k = i \Leftrightarrow x_k \in A_i$ and compute the matrix \hat{P}

$$\hat{P}_{ij} = \frac{C_{ij}}{\sum_j C_{ij}}, \quad C_{ij} = \sum_{k=0}^{N-n} \mathbb{1}_{\{s_k=i\}} \mathbb{1}_{\{s_{k+n}=j\}}. \quad (19)$$

It is well-known [29] that \hat{P} is a maximum likelihood estimator for the full partition MSM transition matrix (4). Similarly one can also compute estimates for a core set MSM by using the definition of milestoning processes [27,28]. That is, if we have core sets C_1, \dots, C_m , a lag time $\tau = nh$ as before, and we define discrete milestoning trajectories by

$$\begin{aligned} s_k^- &= i \Leftrightarrow x_k \in A_i \text{ or came last from } A_i \text{ before time } k \\ s_k^+ &= i \Leftrightarrow x_k \in A_i \text{ or went next to } A_i \text{ after time } k, \end{aligned}$$

we can compute an estimator $\hat{P}_Q = \hat{P}\hat{M}^{-1}$ of the core set MSM matrix (10) by counting transitions:

$$\hat{P}_{ij} = \frac{C_{ij}}{\sum_j C_{ij}}, \quad C_{ij} = \sum_{k=0}^{N-n} \mathbb{1}_{\{s_k^- = i\}} \mathbb{1}_{\{s_{k+n}^+ = j\}}, \quad (20)$$

$$\hat{M}_{ij} = \frac{N_{ij}}{\sum_j N_{ij}}, \quad N_{ij} = \sum_{k=0}^N \mathbb{1}_{\{s_k^- = i\}} \mathbb{1}_{\{s_k^+ = j\}}. \quad (21)$$

187 Since in practice we will only have a finite amount of data available, we will have statistical errors
 188 when constructing an MSM. This is an additional error to the projection error related to the discretization
 189 that we have discussed above. On the other hand, one should note that these errors are *not independent*
 190 of each other. For example, it is clear that if we take a full partition of state space and we let the partition
 191 become arbitrarily fine by letting the number of sets go to infinity, the discretization error will vanish. At
 192 the same time, for a fixed amount of statistics, the statistical error will become arbitrarily large because
 193 we will need to compute more and more estimators for transition events between the increasing number
 194 of sets. For more information on statistical errors we refer to the literature [29,30].

Besides the choice of discretization and the available statistics, the estimates above also depend on a lag time τ . This dependence can be used to validate an MSM by a Chapman Kolmogorov test [29]. This is based on the fact that the MSM matrices approximately form a semi-group for all large enough lag times $\tau > \tau^*$, although for small lag times this is typically not true due to memory effects. These facts also motivate to look at something like an infinitesimal generator that approximately generates these MSM transition matrices for large enough lag times. In [27], two types of generator constructions have been compared for a core set setting. The first generator K is simply constructed from the transition rates between the core sets in the milestoning sense, that is

$$K_{ij} = \lim_{T \rightarrow \infty} \frac{N_{ij}^T}{R_i^T}, \quad i \neq j \quad K_{ii} = - \sum_{j \neq i} K_{ij}, \quad (22)$$

where N_{ij}^T is the amount of time in $[0, T]$ the process has spent on its way from core set C_i to C_j , and R_i^T is the total time in $[0, T]$ the process came last from C_i . On the other hand, one can see [27,31] that

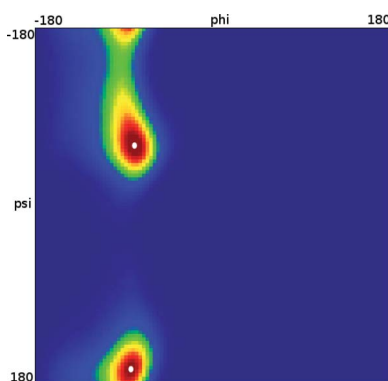
$K^* = KM^{-1}$ with the mass matrix M from above (18) can be interpreted as a projection of the original generator of the process, and also as derivative of the core set MSM from above, i.e.

$$K^* = \lim_{\tau \rightarrow 0} \frac{PM^{-1} - I}{\tau}, \quad (23)$$

195 where P depends on τ (17).

196 Let us now analyze how the choice of core sets, particularly the size of the core sets, influences the
 197 resulting approximation. Therefore, we consider an MD example that was discussed in [27], namely
 198 one molecule of alanine dipeptide monitored via its ϕ and ψ backbone dihedral angles. Two core sets
 199 are defined as balls with radius r around the two points with angular coordinates $x_\alpha = (-80, -60)$ and
 200 $x_\beta = (-80, 170)$. The stationary distribution of the process and the two centers of the core sets x_α, x_β
 201 in the angular space are shown in Fig. 7.

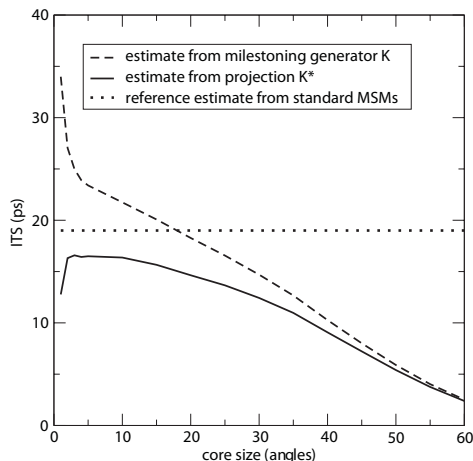
Figure 7. The stationary distribution of alanine dipeptide and the two centers of the core sets x_α, x_β in the angular space as white dots.



202 For computing a reference timescale several MSMs based on full partitions using 10,15, and 250 sets
 203 have been constructed for increasing lag times. In [27], it is shown that in each setting the estimate for
 204 the longest implied timescale of the process converged to ≈ 19 ps for large enough τ . Now the implied
 205 timescales for the two different generators K (22) and K^* (23) are computed. In Fig. 8, the resulting
 206 timescales are plotted against the reference timescale ≈ 19 ps for varying size of the core sets.

207 One can see that the estimate by the milestoning generator K is rather sensitive to the size of core
 208 sets. It overestimates the timescales for small core sizes and underestimates it for larger core sizes. On
 209 the other hand, the projected generator K^* can never overestimate the timescale due to its interpretation
 210 as projection. It is also rather robust against the choice of size of the core sets until the core sets become
 211 too large, e.g. $r > 15$. Then, the discretization becomes close to a full partition discretization using only
 212 two sets. In this case the timescales have to be underestimated heavily because of recrossing phenomena.
 213 On the other hand, the underestimation for very small core sets has to be explained by a lack of statistics.
 214 When the core sets are chosen arbitrarily small, it is clearly more difficult for the process to hit the sets
 215 and therefore transition events become rare. Note that for the straightforward milestoning generator K
 216 the processes seem to become very slow, but for the projected generator $K^* = KM^{-1}$ this effect is
 217 theoretically corrected by the mass matrix M . Nevertheless, in both cases the generation of enough
 218 statistics will be problematic for too small core sets.

Figure 8. Estimate of the implied timescales from K (22), the projected generator K^* (23) and the reference computed from several full partition MSMs.



219 **Further Applications in MD.** Markov State Modelling has been shown to apply successfully to many
 220 different molecular systems like peptides including time-resolved spectroscopic experiments [10–12],
 221 proteins and protein folding [4,9,13], or DNA [32]. In most of the respective publications full partition
 222 MSMs are used and the underlying discretization is based on cluster finding methods, see [29] for a
 223 review. Core set based approaches have been used just recently [10,27].

224 6. MSM for Optimal Control Problems

In this section we will borrow ideas from the previous section and explain how MSMs can be used to discretize optimal control problems that are linear-quadratic in the control variables and which appear in e.g. sampling of rare events. Specifically, we consider the case that $(X_t)_{t \geq 0}$ is the solution of

$$dX_t = (\sqrt{2}u_t - \nabla V(X_t))dt + \sqrt{2\varepsilon}dB_t, \quad (24)$$

with potential V , Brownian motion B_t and temperature $\varepsilon > 0$ as in (16) and an unknown control variable $u: [0, \infty) \rightarrow \mathbb{R}^d$ that is chosen so as to minimize the cost function

$$J(u; x) = \mathbb{E} \left[\int_0^\tau \left(f(X_s) + \frac{1}{2} |u_t|^2 \right) ds \middle| X_0 = x \right]. \quad (25)$$

225 (The factors of $1/2$ and $\sqrt{2}$ in front of the control terms are for notational convenience.) Here $f \geq 0$
 226 is a bounded continuous function called *running cost* and $\tau < \infty$ (a.s.) is a random stopping time that
 227 is determined by X_t hitting a given target set $A \subset E$, i.e. $\tau = \inf\{t > 0: X_t \in A\}$, in other words,
 228 we are interested in controlling $X_t = X_t^u$ until it reaches A . As an example, consider the case $f = 1$
 229 and $A = C_1$ with the potential considered in Figure 5, which amounts to the situation that one seeks to
 230 minimize the time to reach the core set C_1 by tilting the potential towards the target set C_1 ; tilting the
 231 potential too much is prevented by the quadratic penalization term in the cost functional that grows when
 232 too much force is applied.

233 Other choices of f in (24) result in alternative applications. One obvious application would be to
 234 set $\tau = T$ to a fixed time and f to the characteristic function of the complement of a conformation set
 235 C , $f = \mathbb{1}_{E \setminus C}$. In this case, minimization of J wrt. the control u_t would mean maximization of the
 236 probability to find the system in the conformation C until time T under a penalty on the external work
 237 done to the system. See [14] for more details on such applications.

238 There are other types of cost functions J one might consider, e.g. control until a deterministic finite
 239 time $\tau = T$ is reached, or even $\tau \rightarrow \infty$, and the construction would follow analogously. For compactness
 240 we consider here only cost functions as in (25).

Optimal control and equilibrium expectation values. It turns out that when minimizing J it is sufficient to consider control strategies that are Markovian and depend only on X_t , i.e. we consider feedback laws of the form $u_t = \alpha(X_t)$ for some smooth function $\alpha: E \rightarrow \mathbb{R}^d$. Moreover only controls with finite energy are considered, for otherwise $J(u; x) = \infty$. For control problems of the form (24)–(25) the optimal feedback function can be shown to be $\alpha^*(x) = -\sqrt{2}\nabla W$ where W is the value function or optimal-cost-to-go [1,15]

$$W(x) = \min_u J(u; x) \quad (26)$$

with the minimum running over all admissible Markovian feedback strategies. It can be shown that W satisfies the following dynamic programming equation of Hamilton-Jacobi-Bellman type (see [33]):

$$\begin{aligned} LW(x) - |\nabla W(x)|^2 + f &= 0 \\ W|_A &= 0, \end{aligned} \quad (27)$$

with the second-order differential operator

$$L = \varepsilon\Delta - \nabla V \cdot \nabla$$

that is the infinitesimal generator of the process X_t for $u = 0$. If the value function W is known, it can be plugged into the equation of motion which then turns out to be of the form

$$dX_t^* = -\nabla U(X_t^*)dt + \sqrt{2\varepsilon}dB_t, \quad (28)$$

with the new potential

$$U(x) = V(x) + 2W(x).$$

The difficulty is that equation (27) is a nonlinear partial differential equation and for realistic high-dimensional systems it is not at all obvious how to discretize it, employing any kind of state space partitioning. It has been demonstrated in [14,15] that (27) can be transformed into a linear equation by a logarithmic transformation. Setting $W(x) = -\varepsilon \log \phi(x)$ it readily follows, using chain rule and equation (27), that ϕ solves the linear equation

$$\begin{aligned} (L - \varepsilon^{-1}f)\phi &= 0 \\ \phi|_A &= 1. \end{aligned} \quad (29)$$

The last equation is linear and can be solved by using MSMs as we will show below. Moreover, by the Feynman-Kac theorem [34], the solution to (29) can be expressed as

$$\phi(x) = \mathbb{E} \left[\exp \left(-\frac{1}{\varepsilon} \int_0^\tau f(X_t) dt \right) \middle| X_0 = x \right], \quad (30)$$

where X_t solves the control-free equation

$$dX_t = -\nabla V(X_t)dt + \sqrt{2\varepsilon}dB_t.$$

241 That is, the optimal control for (24) can be computed by solving (29) which can be done *in principle* via
 242 Monte-Carlo approximation of the expected value in (30) if critical slowing down by rare events can be
 243 avoided.

Remark. The optimization problem (26) admits an interpretation in terms of entropy minimization: Let $Q = Q_x^u$ and $P = Q_x^0$ denote the path probability measures of controlled and uncontrolled trajectories starting at x at time $t = 0$, and set

$$Z = \int_0^\tau f(X_s) ds,$$

then it follows that we can write

$$W(x) = \min_{Q \ll P} J(u; x), \quad J(u; x) = \int \left\{ Z + \varepsilon \log \left(\frac{dQ}{dP} \right) \right\} dQ, \quad (31)$$

where the notation “ $Q \ll P$ ” means that Q has a density¹ with respect to P . It turns out that for every such Q there is exactly one control strategy u such that $Q = Q_x^u$ is generated by (24), in this sense the notation in (31) is meaningful. The second term

$$H(Q||P) = \varepsilon \int \log \left(\frac{dQ}{dP} \right) dQ$$

244 is the relative entropy or Kullback-Leibler divergence between Q and P . For details on this matter that
 245 are based on Girsanov transformations for stochastic differential equations we refer to [35] or the article
 246 [1] in this special issue.

247 7. MSM Discretization of Optimal Control Problems

248 The basic idea is now to choose a subspace $D \subset L^2(\mu)$ with basis χ_1, \dots, χ_n as in Markov state
 249 modelling and then discretize the dynamic programming equation (27) of our optimal control problem by
 250 projecting the equivalent log transformed equation (29) onto that subspace. As we will see the resulting
 251 discrete matrix equation can be transformed back into an optimal control problem for a discrete Markov
 252 jump process (MJP).

253 We will do this construction for the full partition case $\chi_i = \mathbb{1}_{A_i}$ and the core set case $\chi_i = q_i$ discussed
 254 earlier. We will see that in both cases, we arrive at a structure-preserving discretization of the original
 255 optimal control problem where the states of the corresponding MJP will be related to the partition subsets
 256 A_i . The first case will give us back a well-known lattice discretization for continuous control problems,
 257 the Markov chain approximation [36]. This is illustrated in the following diagram:

¹That is, the density function dQ/dP exists, is almost everywhere positive and normalized.

$$\begin{array}{ccc}
& \text{SDE} & \text{MJP} \\
\text{Linear equation} & L\phi = \epsilon^{-1}f\phi & \xrightarrow[D \subset V]{\text{discretize}} G\hat{\phi} = \epsilon^{-1}\hat{f}\hat{\phi} \\
& \updownarrow W = -\epsilon \log \phi & \updownarrow \hat{W} = -\epsilon \log \hat{\phi} \\
\text{Control Problem} & W = \min_u J(u) & \xrightarrow{\text{?}} \hat{W} = \min_v \hat{J}(v)
\end{array}$$

258 **Subspace projection.** The key steps for the discretization is that we pick a suitable subspace $D \subset$
259 $L^2(\mu)$ that is adapted to the boundary value problem (29). Specifically, we require that the subspace
260 contains the constant function $\mathbb{1} \in D$ and that it gives a good representation of the most dominant
261 metastable sets. To this end we choose basis functions $\chi_1, \dots, \chi_{n+1}$ with the following properties:

262 (S1) The χ_i form a partition of unity, that is $\sum_{i=1}^{n+1} \chi_i = \mathbb{1}$.

263 (S2) The χ_i are adapted to the boundary conditions in (29), that is $\chi_{n+1}|_A = 1$ and $\chi_i|_A = 0$ for
264 $i \in \{1, \dots, n\}$.

Now let Q be the orthogonal projection onto D , and define the matrices

$$F_{ij} = \frac{\langle \chi_i, f\chi_j \rangle}{\langle \chi_i, \mathbb{1} \rangle}, \quad K_{ij} = \frac{\langle \chi_i, L\chi_j \rangle}{\langle \chi_i, \mathbb{1} \rangle}.$$

Now, if ϕ solves the linear boundary value problem (29), then the coefficients $\hat{\phi}_1, \dots, \hat{\phi}_{n+1}$ of its finite-dimensional representation $Q\phi = \sum_j \hat{\phi}_j \chi_j$ on the subspace D satisfy the constrained linear system

$$\begin{aligned}
\sum_{j=1}^{n+1} (K_{ij} - \epsilon^{-1}F_{ij}) \hat{\phi}_j &= 0, \quad i \in \{1, \dots, n\} \\
\hat{\phi}_{n+1} &= 1,
\end{aligned} \tag{32}$$

that is the discrete analogue of (29). The discrete solution $\hat{\phi} = Q\phi$ is optimal in the sense of being the best approximation of ϕ in the energy norm, i.e.,

$$\|\phi - \hat{\phi}\|_A = \inf_{\psi \in D} \|\phi - \psi\|_A, \tag{33}$$

where

$$\|\phi\|_A^2 = \langle \phi, (\epsilon^{-1}f - L)\phi \rangle$$

is the energy norm on $L^2(\mu)$, and the infimum runs over all functions $\psi \in L^2(\mu)$ that are of the form $\psi(x) = \sum_j \psi_j \chi_j(x)$ with coefficients $\psi_j \in \mathbb{R}$. This is a standard result about projections of PDEs, see [37] for details.² In analogy with equation (14) we can use the above result to get the error estimate

$$\|\phi - \hat{\phi}\|_\mu^2 \leq \left(1 + \frac{1}{\delta^2} \|QAQ^\perp\|^2\right) \inf_{\psi \in D} \|\phi - \psi\|_\mu^2 \tag{34}$$

²By the same argument as in the previous sections $A = \epsilon^{-1}f - L$ is symmetric and positive definite as an operator on the weighted Hilbert space $L^2(\mu)$. Moreover $\|\phi\|_A^2 = \epsilon^{-1}\langle \phi, f\phi \rangle + \epsilon\langle \nabla\phi, \nabla\phi \rangle$.

265 where $A = \varepsilon^{-1}f - L$ is a shorthand for the operator appearing in (29) and the constant $\delta > 0$ is defined
 266 such that $\|v\|_A^2 \geq \delta \|v\|_\mu^2$ holds for all $v \in L^2(\mu)$; see [38]. The bottom line of (33) is that discretizing (29)
 267 via (32) minimizes the projection error measured in the energy norm. Since all functions are μ -weighted,
 268 the approximation will be good in regions visited with high probability and less good in regions with
 269 lower probability. The error estimate (34) is along the lines of the MSM approximation result: If we
 270 switch to the norm on $L^2(\mu)$, the function $\hat{\phi} = Q\phi$ is still almost the best approximation of ϕ , provided
 271 that A leaves the subspace D almost invariant. As was pointed out earlier this is exactly the case when
 272 the χ_i are close to the eigenfunctions of A (e.g., when the system is metastable).

Properties of the projected problem. We introduce now the diagonal matrix Λ with entries $\Lambda_{ii} = \sum_j F_{ij}$ (zero otherwise) and the full matrix $G = K - \varepsilon^{-1}(F - \Lambda)$, and rearrange (32) as follows:

$$\sum_{j=1}^{n+1} (G_{ij} - \varepsilon^{-1}\Lambda_{ij}) \hat{\phi}_j = 0, \quad i \in \{1, \dots, n\} \quad (35)$$

$$\hat{\phi}_{n+1} = 1,$$

273 This equation can be given a stochastic interpretation. To this end let us introduce the vector $\pi \in \mathbb{R}^{n+1}$
 274 with nonnegative entries $\pi_i = \langle \chi_i, \mathbf{1} \rangle$ and notice that $\sum_i \pi_i = 1$ follows immediately from the fact that
 275 the basis functions χ_i form a partition of unity, i.e. $\sum_i \chi_i = \mathbf{1}$. This implies that π is a probability
 276 distribution on the discrete state space $\hat{E} = \{1, \dots, n+1\}$. We summarise properties of the matrices K ,
 277 F and G , see also [38]:

(M1) K is a generator matrix of a MJP $(\hat{X}_t)_{t \geq 0}$ (i.e., K is a real-valued square matrix with row sum zero and positive off-diagonal entries) with stationary distribution π that satisfies detailed balance

$$\pi_i K_{ij} = \pi_j K_{ji}, \quad i, j \in \hat{E}$$

278 (M2) $F \geq 0$ (entry-wise) with $\pi_i F_{ij} = \pi_j F_{ji}$ for all $i, j \in \hat{E}$.

279 (M3) G has row sum zero and satisfies $\pi^T G = 0$ and $\pi_i G_{ij} = \pi_j G_{ji}$ for all $i, j \in \hat{E}$; furthermore there
 280 exists a constant $0 < C < \infty$ such that $G_{ij} \geq 0$ for all $i \neq j$ if $\|f\|_\infty \leq C$. In this case equation
 281 (35) admits a unique and strictly positive solution $\hat{\phi} > 0$.

It follows that if the running costs f are such that (M3) holds, then G is a generator matrix of a MJP that we shall denote by $(\hat{X}_t)_{t \geq 0}$, and (35) has a unique and positive solution. In this case the logarithmic transformation $\hat{W} = -\varepsilon \log \hat{\phi}$ is well-defined. It was shown in [39] that \hat{W} can be interpreted as the value function of a Markov decision problem with cost functional (cf. also [33])

$$\hat{J}(v; i) = \mathbb{E} \left[\int_0^\tau \left(f(\hat{X}_s) + k(\hat{X}_s, v_s) \right) ds \middle| \hat{X}_0 = i \right] \quad (36)$$

that is minimized over the set of Markovian control strategies $v: \hat{E} \rightarrow (0, \infty)$ subject to the constraint that the controlled process $\hat{X}_t = \hat{X}_t^v$ is generated by G^v where

$$G_{ij}^v = \begin{cases} v(i)^{-1} G_{ij}, & i \neq j \\ -\sum_{j \neq i} G_{ij}^v, & i = j \end{cases} \quad (37)$$

with stopping time $\tau = \inf\{t > 0: \hat{X}_t = n + 1\}$ and running costs

$$\hat{f}(i) = \Lambda_{ii}, \quad k(i, v) = \varepsilon \sum_{j \neq i} G_{ij} \left\{ \frac{v(j)}{v(i)} \left[\log \frac{v(j)}{v(i)} - 1 \right] + 1 \right\}. \quad (38)$$

Properties of the projected problem, cont'd. From [39] we know that the optimal cost

$$\hat{W}(i) = \min_v \hat{J}(v; i)$$

is given by $\hat{W} = -\varepsilon \log \hat{\phi}$ where $\hat{\phi}$ solves (35), with the optimal feedback strategy given by $v^*(i) = \hat{\phi}_i$ (see [33]). We list additional properties:

(i) The v -controlled system has the unique invariant distribution

$$\pi^v = (\pi_1^v, \dots, \pi_{n+1}^v), \quad \pi_i^v = \frac{v(i)^2 \pi_i}{Z_v}$$

with Z_v an appropriate normalization constant; in terms of the value function $\pi^* = \pi^{v^*}$ reads

$$\pi^* = (\pi_1^*, \dots, \pi_{n+1}^*), \quad \pi_i^* = \frac{1}{Z_*} e^{-2\varepsilon^{-1}\hat{W}(i)} \pi_i.$$

(ii) G^v is reversible and stationary with respect to π^v , i.e., $\pi_i^v G_{ij}^v = \pi_j^v G_{ji}^v$ for all $i, j \in \hat{E}$.

(iii) \hat{J} admits the same interpretation as (31) in terms of the relative entropy:

$$\hat{W}(i) = \min_{Q \ll P} \hat{J}(v; i), \quad \hat{J}(v; i) = \int \left\{ \hat{Z} + \varepsilon \log \left(\frac{dQ}{dP} \right) \right\} dQ$$

where P denotes expectation with respect to the uncontrolled MJP \hat{X}_t starting at $\hat{X}_0 = i$, Q denotes the path measure of the corresponding controlled process with generator G^v and

$$\hat{Z} = \int_0^\tau \hat{f}(\hat{X}_s) ds.$$

A few remarks seem in order: Item (i) of the above list is in accordance with the continuous setting, in which the optimally controlled dynamics is governed by the new potential $U = V + 2W$ and has the stationary distribution $\mu^* \propto \exp(-2\varepsilon^{-1}W)\mu$ with μ being the stationary distribution of the uncontrolled process. Hence the effect of the control on the invariant distribution is the same in both cases. Further note that optimal strategies change the jump rates according to

$$G_{ij}^{v^*} = G_{ij} e^{-\varepsilon^{-1}(\hat{W}(j) - \hat{W}(i))}, \quad (39)$$

that is \hat{W} acts as an effective potential as in the continuous case, and the change in the jump rates can be interpreted in terms of Kramer's law for this effective potential.

This completes our derivation of the discretized optimal control problem, and we now compare it with the continuous problem we started with for the case of a full partition of E and a core set partition of E .

289

290

291 **8. Markov Chain Approximations and Beyond**

292 **Full partitions.** Let E be fully partitioned into disjoint sets A_1, \dots, A_{n+1} with centers x_1, \dots, x_{n+1}
 293 and such that $A_{n+1} := A$, and define $\chi_i := \chi_{A_i}$. These χ_i satisfy the assumptions (S1) and (S2) discussed
 294 in section 7. Since they are not overlapping, F is diagonal, and

$$\hat{f}(i) = \frac{1}{\pi_i} \int_{A_i} f(x) \mu(x) dx = \mathbb{E}_\mu[f(X_t) | X_t \in A_i] \quad (40)$$

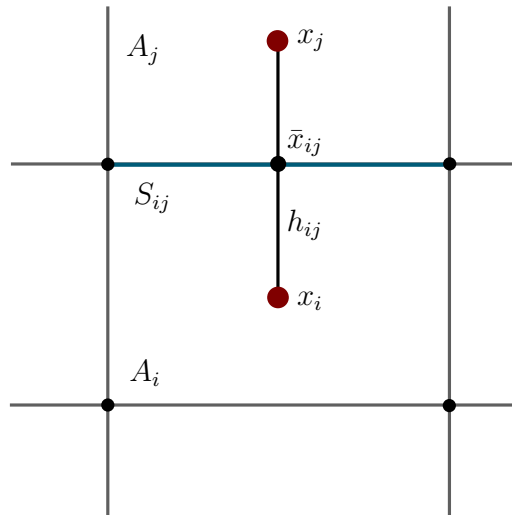
295 is just obtained by averaging $f(x)$ over the cell A_i . (40) is also a sampling formula for $\hat{f}(i)$. It follows
 296 directly that $G = K$, and in particular (M3) holds for any f . One can show that K has components

$$K_{ij} \approx \frac{1}{\Delta_{ij}} e^{-\beta(V(\bar{x}_{ij}) - V(x_i))}, \quad \Delta_{ij}^{-1} = \beta^{-1} \frac{m(S_{ij})}{m(h_{ij})m(A_i)} \quad (41)$$

297 if i and j are neighbours ($K_{ij} = 0$ otherwise). Here m is the Lebesgue measure, and h_{ij} , S_{ij} and \bar{x}_{ij}
 298 are defined as in figure 9. K is the generator of a MJP on the cells A_i and coincides with the so-called
 299 *finite volume approximation* of L discussed in [40]. It is reversible with stationary distribution

$$\pi_i = \int_{A_i} d\mu \approx m(A_i) e^{-\beta V(x_i)}.$$

Figure 9. The mesh for the full partition.



300 One can show that the approximation error vanishes for $n \rightarrow \infty$. K and π can be computed from the
 301 potential V and the geometry of the mesh. By inspecting (12) and (13), we see that K is connected to
 302 the transition matrix P^τ of a full partition MSM with lagtime τ by

$$\lim_{\tau \rightarrow 0} \frac{1}{\tau} (P_{ij}^\tau - M_{ij}) = \lim_{\tau \rightarrow 0} \frac{1}{\pi_i} \langle \chi_i, \frac{1}{\tau} (T_\tau - \mathbb{1}) \chi_j \rangle = \frac{1}{\pi_i} \langle \chi_i, L \chi_j \rangle = K_{ij},$$

303 thus K is the generator of the semigroup of transition matrices P^τ . Therefore we could obtain K by
 304 sampling in the same way we obtained P^τ through equation (19) in section 5. This is difficult however

305 due to recrossing problems for small τ , see e.g. [41]. Finally, let us note in passing that we can drastically
 306 simplify k^v if the cells A_i are boxes of length h . Denote the elementary lattice vectors by e_n . Then

$$k^v(i) = \frac{1}{2}|u^v(i)|^2 + \mathcal{O}(h), \quad u_n^v(i) := \frac{1}{\sqrt{2}} \frac{\epsilon}{2h} (\log v(i + e_n) - \log v(i - e_n))$$

307 which establishes the connection to the continuous case. But more is true: The whole discrete control
 308 problem reduces to first order in h to the well-known Markov chain approximation (MCA) [36], which
 309 allows us to use convergence theory for MCAs to conclude that for $n \rightarrow \infty$, optimal control and value
 310 function of the discrete control problem converge to their continuous counterparts. More details can be
 311 found in [38].

312 **Core set partition.** Now we choose core sets C_1, \dots, C_{n+1} with $C_{n+1} = A$ and we let $\chi_i = q_i$ to be
 313 the committor function of the process with respect to C_i as in section 4. These χ_i satisfy the assumptions
 314 (S1) and (S2) discussed in section 7. Recall the definition of the forward and backward milestoning
 315 process \tilde{X}_t^\pm from (18). The discrete costs can be written as

$$\hat{f}(i) = \frac{1}{\pi_i} \langle q_i, f \sum_j q_j \rangle = \int \nu_i(x) f(x) dx = \mathbb{E}_\mu \left[f(X_t) \middle| \tilde{X}_t^- = i \right] \quad (42)$$

316 where $\nu_i(x) = \frac{q_i(x)\mu(x)}{\pi_i} = \mathbb{P}(X_t = x | \tilde{X}_t^- = i)$ is the probability density of finding the system in state
 317 x given that it came last from i . Hence $\hat{f}(i)$ is the average costs conditioned on the information $\tilde{X}_t^- = i$,
 318 i.e. X_t came last from A_i , which is the natural extension to the full partition case where $\hat{f}(i)$ was the
 319 average costs conditioned on the information that $X_t \in A_i$.

320 The matrix $K = \pi_i^{-1} \langle q_i, Lq_j \rangle$ is reversible with stationary distribution

$$\pi_i = \langle q_i, \mathbb{1} \rangle = \mathbb{P}_\mu(\tilde{X}_t^- = i)$$

321 and is related to core MSMs again:

$$K = \lim_{\tau \rightarrow 0} \frac{1}{\tau} (P^\tau - M)$$

322 where P^τ and M are now the matrices for core MSMs as in (18). Formally, K is the generator of the
 323 P^τ , but these do not form a semigroup since $M \neq \mathbb{1}$, and therefore we cannot interpret K directly as
 324 e.g. the generator of \tilde{X}_t^- . Nevertheless, the entries of K are the transition rates between the core sets as
 325 defined in transition path theory [42]. We can sample P^τ and M using (20) and (21), and because we
 326 used an incomplete partition, the recrossing problem is removed, and there is no difficulty in sampling
 327 P^τ for all lagtimes τ and therefore K directly. It is worth noting that F can also be sampled:

$$F_{ij} = \mathbb{E}_\mu \left[f(X_t) \chi_{\{\tilde{X}_t^+ = j\}} \middle| \tilde{X}_t^- = i \right]$$

328 Therefore, as in the construction of core MSMs, we do not need to compute committor functions
 329 explicitly. Note however that $G \neq L$, there is a reweighting due to the overlap of the q_i 's which causes
 330 F to be nondiagonal. This reweighting is the surprising bit of this discretization. From properties
 331 (M1)-(M3) from section 7 we see however that G and K are both reversible with stationary distribution
 332 π . Finally, note that if the cost function $f(x)$ doesn't satisfy $\|f\|_\infty \leq C$ from (M3), G will not even be

333 a generator matrix. In this case (32) still has a solution $\hat{\phi}$ which is the best approximation to ϕ , but this
 334 solution may not be unique, it may not satisfy $\hat{\phi} > 0$, and we have no interpretation as a discrete control
 335 problem.

336 9. Numerical Results

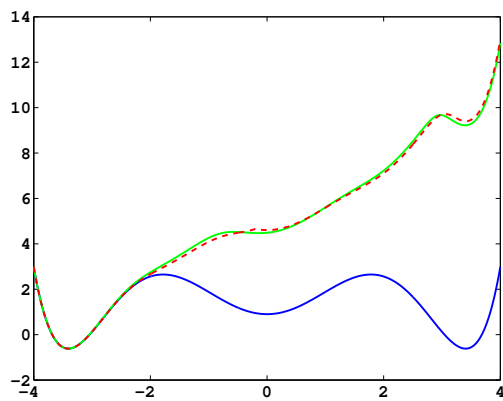
337 9.1. 1D Potential Revisited

338 Firstly, we study diffusion in the triple well potential which is presented in Figure 2. This potential
 339 has three minima at approximately $x_{0/1} = \pm 3.4$ and $x_2 = 0$. We choose the three core sets $C_i =$
 340 $[x_i - \delta, x_i + \delta]$ around the minima with $\delta = 0.2$. Take τ to be the first hitting time of C_0 . We are
 341 interested in the moment generating function $\phi(x) = \mathbb{E} \left[e^{-\epsilon^{-1} \sigma \tau} \right]$ of passages into C_0 and the cumulant
 342 generating function $W = \epsilon \log \phi$. This is of the form (30) for $A = C_i$ and $f = \sigma$ a constant function.

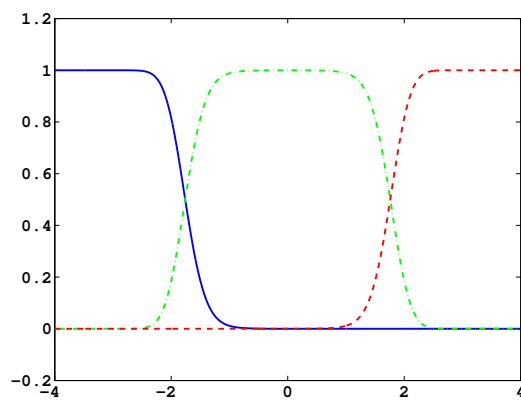
343 In figure 10a the potential V and effective potential U are shown for $\beta = 2$ and $\sigma = 0.08$ (solid
 344 lines), cf. equation (28). One can observe that the optimal control effectively lifts the second and third
 345 well up which means that the optimal control will drive the system into C_0 very quickly. The reference
 346 computations here have been carried out using a full partition FEM discretization of (29) with a lattice
 347 spacing of $h = 0.01$. Now we study the MJP approximation constructed via the committor functions
 348 shown in Figure 10b. These span a three-dimensional subspace, but due to the boundary conditions
 349 the subspace D of the method is actually two-dimensional. The dashed line in Figure 10a gives the
 350 approximation to U calculated by solving (35). We can observe extremely good approximation quality,
 351 even in the transition region. In Figure 10c the approximation to the optimal control $\alpha^*(x)$ (solid line)
 352 and its approximation $\hat{\alpha}^* = -\sqrt{2} \nabla \hat{W}$ (dashed line) are shown. The core sets are shown in blue. We can
 353 observe jumps in $\hat{\alpha}^*$ at the left boundaries of the core sets. This is to be expected and comes from the
 354 fact that the committor functions are not smooth at the boundaries of the core sets, but only continuous.
 355 Therefore the approximation to U is continuous, but the approximation to α^* is not.

356 Next we construct a core MSM to sample the matrices K and F . 100 trajectories of length $T = 20000$
 357 were used to build the MSM. In Figure 10d, W and its estimate using the core MSM is shown for $\epsilon = 0.5$
 358 and different values of σ . Each of the 100 trajectories has seen about four transitions. For comparison,
 359 a direct sampling estimate of W using the same data is shown (green). The direct sampling estimate
 360 suffers from a large bias and variance and is practically useless. In contrast, the MSM estimator for W
 361 performs well for all considered values of σ and always its variance is significantly small. The constant
 362 C which ensures $\hat{\phi} > 0$ when $\sigma \leq C$ is approximately 0.2 in this case. This seems restrictive but still
 363 allows to capture all interesting information about ϕ and W .

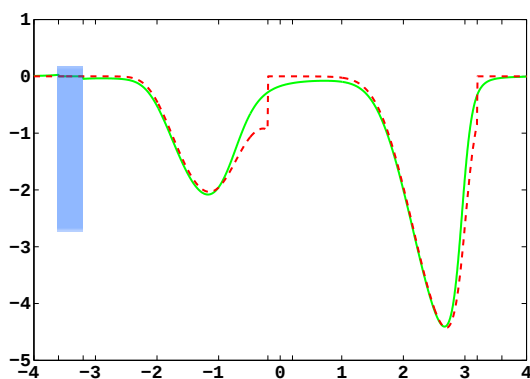
Figure 10. Three well potential example for $\epsilon = 0.5$ and $\sigma = 0.08$. (a) Potential $V(x)$ (blue), effective potential $U = V + 2W$ (green) and approximation of U with committors (dashed red). (b) The three committors $q_1(x)$, $q_2(x)$ and $q_3(x)$. (c) The optimal control $\alpha^*(x)$ (solid line) and its approximation (dashed line). Core sets are shown in blue. (d) Optimal cost W for $\beta = 2$ as a function of σ . Blue: Exact solution. Red: Core MSM estimate. Green: Direct sampling estimate.



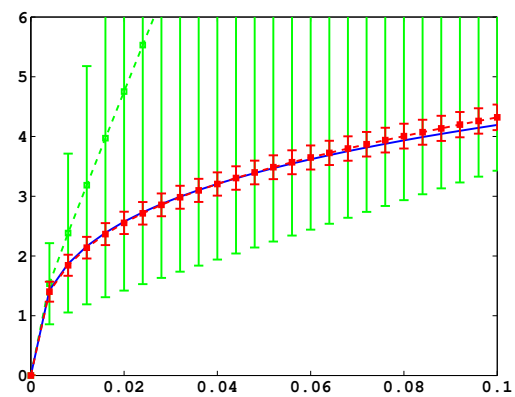
(a)



(b)



(c)



(d)

364 *9.2. Alanine Dipeptide*

365 Lastly, we study α - β -transitions in Alanine dipeptide, a well-studied test system for Molecular
 366 Dynamics applications. We use a $1\mu s$ long trajectory simulated with the CHARMM 27 force field. The
 367 conformational dynamics is monitored as usual via the backbone dihedral angles ϕ and ψ . The data was
 368 first presented in [43]. We construct a full partition MSM with 250 clusters using k -means clustering.
 369 We are interested in the MFPT $\hat{t}(i) = \mathbb{E}_i[\tau_\alpha]$ where τ_α is the first hitting time of the α conformation,
 370 which we define as a circle with radius $r = 45$ around $(\phi_\alpha, \psi_\alpha) = (-80, -60)$. The MFPT vector \hat{t}
 371 solves the boundary value problem

$$K\hat{t} = -1 \text{ outside of } \alpha, \quad \hat{t} = 0 \text{ in } \alpha,$$

372 but since K is not available directly via sampling, we have to consider the equation

$$\frac{1}{\tau}(P^\tau - 1)\hat{t} = -1 \text{ outside of } \alpha, \quad \hat{t} = 0 \text{ in } \alpha$$

373 instead. The result will depend on the choice of lagtime τ . In Figure 11a, the results are shown for
 374 $\tau = 5$, we can identify the β -structure as the red cloud of clusters where $\hat{t}(i)$ is approximately constant.
 375 In 11b, $\hat{t}_{\beta\alpha} = \mathbb{E}(\hat{t}(i)|i \in \beta)$ is shown as a function of τ . We observe a linear behavior for large τ which
 376 is due to the linear error introduced in the replacement of K with $\frac{1}{\tau}(P^\tau - 1)$ and a nonlinear drop for
 377 small τ which is due to Non-Markovianity. Our best guess is therefore a linear interpolation to $\tau = 0$,
 378 which is indicated by the solid line. The result is $\hat{t}_{\beta\alpha}^{(0)} = 35.5ps$. As a comparison the reference value
 379 $\hat{t}_{\beta\alpha}^{ref} = 36.1ps$ from [43] is shown as a dashed line. It was computed in [43] as an inverse rate, using
 380 the slowest ITS and information about the equilibrium weights of the α and β structure. We see very
 381 good agreement. The result is of course dependent though on the assignment of clusters to the α and β
 382 structure. Some tests show that $\hat{t}_{\beta\alpha}^{(0)}$ as computed with the interpolation method is fairly insensitive to this
 383 choice.

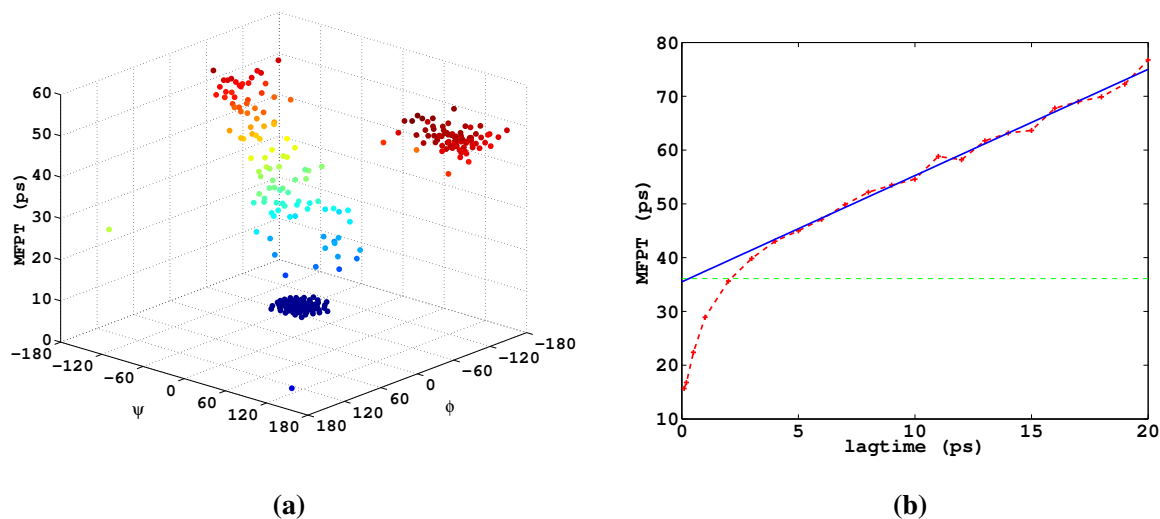
384 In [14] it is demonstrated how to use the method presented herein for maximizing the population of
 385 the α -conformation of Alanin dipeptide based on the MSM used here.

386 **10. Conclusion**

387 In this article, we have discussed an approach to overcome direct sampling issues of rare events
 388 in molecular dynamics based on spatial discretization of the molecular state space. The strategy is to
 389 define a discretization by subsets of state space such that the sampling effort with respect to transitions
 390 between the sets is much lower than the direct estimation of the rare events under consideration. That
 391 is, without having to simulate rare events we construct a so called Markov State Model, a Markov chain
 392 approximation to the original dynamics. Since the state space of the MSM is finite, we can then calculate
 393 the properties of interest by simply solving linear systems of equations. Of course, it is crucial that these
 394 properties of the MSM can be related to the rare event properties of the original process that we have not
 395 been able to sample directly.

396 This is why we have analyzed the approximation quality of MSMs in the first part of the article. We
 397 have used the interpretation of MSMs as projections of the transfer operator to (1) derive conditions that

Figure 11. Dipeptide example. (a) MFPT from β to α in ϕ - ψ space for $\tau = 5$. The red cloud to the right is the β -structure. (b) MFPT as a function of τ (dashed line) and linear interpolation to $\tau = 0$ (solid line). Green dashed line: Reference computed via slowest ITS.



398 guarantee an accurate reproduction of the dynamics, and (2) show how to construct models based on a
399 core set discretization by leaving the state space partly undiscretized.

400 In the second part of the article, we have used the concept of MSM discretization to solve MD optimal
401 control problems in which one computes the optimal external force that drives the molecular system
402 to show an optimized behavior (maximal possible population in a conformation; minimal mean first
403 passage time to a certain conformation) under certain constraints. We have demonstrated that the spatial
404 discretization underlying an MSM turns the high-dimensional continuous optimal control problem into a
405 rather low-dimensional discrete optimal control problem of the same form that can be solved efficiently.
406 This result allows two different types of application: (1) If one can construct an MSM for a molecular
407 system in equilibrium, then one can use it to compute optimal controls that extremize a given costs
408 criterion. (2) If an MSM can be computed based on transition probabilities between neighboring core
409 sets alone then the rare event statistics for transitions between strongly separated metastable states
410 of the system can be computed from an associated optimal control problem that can be solved after
411 discretization using the pre-computed MSM.

412 References

- 413 1. Hartmann, C.; Banisch, R.; Sarich, M.; Badowski, T.; Schütte, C. Characterization of Rare
414 Events in Molecular Dynamics. *Entropy* **2013**. submitted.
- 415 2. Metzner, P.; Schütte, C.; Vanden-Eijnden, E. Transition Path Theory for Markov Jump Processes.
416 *Multiscale Modeling and Simulation* **2009**, 7, 1192–1219.
- 417 3. Pan, A.C.; Roux, B. Building Markov state models along pathways to determine free energies
418 and rates of transitions. *J. Chem. Phys.* **2008**, 129, 064107+.

- 419 4. Noé, F.; Schütte, C.; Vanden-Eijnden, E.; Reich, L.; Weikl, T. Constructing the full ensemble of
420 folding pathways from short off-equilibrium trajectories. *PNAS* **2009**, *106*(45), 19011–19016.
- 421 5. Schütte, C. Conformational Dynamics: Modelling, Theory, Algorithm, and Applications to
422 Biomolecules. Habilitation thesis, Fachbereich Mathematik und Informatik, FU Berlin, 1998.
- 423 6. Schütte, C.; Fischer, A.; Huisinga, W.; Deuffhard, P. A Direct Approach to Conformational
424 Dynamics based on Hybrid Monte Carlo. *J. Comput. Phys.* **1999**, *151*, 146–168.
- 425 7. Prinz, J.; Wu, H.; Sarich, M.; Keller, B.; Fischbach, M.; Held, M.; Schuette, C.; Chodera, J.D.;
426 Noe, F. Markov models of molecular kinetics: Generation and Validation. *J. Chem. Phys.* **2010**.
427 submitted.
- 428 8. Noé, F.; Horenko, I.; Schütte, C.; Smith, J.C. Hierarchical Analysis of Conformational Dynamics
429 in Biomolecules: Transition Networks of Metastable States. *J. Chem. Phys.* **2007**, *126*, 155102.
- 430 9. Chodera, J.D.; Dill, K.A.; Singhal, N.; Pande, V.S.; Swope, W.C.; Pitera, J.W. Automatic
431 discovery of metastable states for the construction of Markov models of macromolecular
432 conformational dynamics. *J. Chem. Phys.* **2007**, *126*, 155101.
- 433 10. Buchete, N.V.; Hummer, G. Coarse Master Equations for Peptide Folding Dynamics. *J. Phys.*
434 *Chem. B* **2008**, *112*, 6057–6069.
- 435 11. Prinz, J.H.; Keller, B.; Noé, F. Probing molecular kinetics with Markov models: Metastable
436 states, transition pathways and spectroscopic observables. *Phys. Chem. Chem. Phys.* **2011**,
437 *13*, 16912–16927.
- 438 12. Keller, B.; Prinz, J.H.; Noé, F. Markov models and dynamical fingerprints: Unraveling the
439 complexity of molecular kinetics. *Chem. Phys. (in press)* **2011**. Accepted 31. August 2011.
- 440 13. Bowman, G.; Volez, V.; Pande, V.S. Taming the complexity of protein folding. *Current Opinion*
441 *in Structural Biology* **2011**, *21*, 4–11.
- 442 14. Schütte, C.; Winkelmann, S.; Hartmann, C. Optimal control of molecular dynamics using Markov
443 state models. *Math. Program. Series B* **2012**, *134*, 259–282.
- 444 15. Hartmann, C.; Schütte, C. Efficient rare event simulation by optimal nonequilibrium forcing. *J.*
445 *Stat. Mech. Theor. Exp.* **2012**, p. 11004.
- 446 16. Jäger, M.; Zhang, Y.; Bieschke, J.; Nguyen, H.; Dendle, M.; Bowman, M.E.; Noel, J.P.; Gruebele,
447 M.; Kelly, J.W. Structure-function-folding relationship in a WW domain. *Proc. Natl. Acad. Sci.*
448 *USA* **2006**, *103*, 10648–10653.
- 449 17. Kobitski, A.Y.; Nierth, A.; Helm, M.; Jäschke, A.; Nienhaus, G.U. Mg²⁺ dependent folding of
450 a Diels-Alderase ribozyme probed by single-molecule FRET analysis. *Nucleic Acids Res.* **2007**,
451 *35*, 2047–2059.
- 452 18. Fischer, S.; Windshuegel, B.; Horak, D.; Holmes, K.C.; Smith, J.C. Structural mechanism of the
453 recovery stroke in the Myosin molecular motor. *Proc. Natl. Acad. Sci. USA* **2005**, *102*, 6873–
454 6878.
- 455 19. Noé, F.; Krachtus, D.; Smith, J.C.; Fischer, S. Transition Networks for the Comprehensive
456 Characterization of Complex Conformational Change in Proteins. *J. Chem. Theo. Comp.* **2006**,
457 *2*, 840–857.
- 458 20. Ostermann, A.; Waschipky, R.; Parak, F.G.; Nienhaus, U.G. Ligand binding and conformational
459 motions in myoglobin. *Nature* **2000**, *404*, 205–208.

- 460 21. Huisinga, W. Metastability of Markovian Systems A transfer operator based approach in
461 application to molecular dynamics. Phd thesis, Fachbereich Mathematik und Informatik, FU
462 Berlin, 2001.
- 463 22. Bovier, A.; Eckhoff, M.; Gaynard, V.; Klein, M. Metastability in reversible diffusion processes.
464 I. Sharp asymptotics for capacities and exit times. *J. Eur. Math. Soc. (JEMS)* **2004**, *6*, 399–424.
- 465 23. Voronoi, M.G. Nouvelles applications des parametres continus a la theorie des formes
466 quadratiques. *J. Reine Angew. Math.*, *134*, 198–287.
- 467 24. Sarich, M.; Noé, F.; Schütte, C. On the Approximation Quality of Markov State Models.
468 *Multiscale Modeling and Simulation* **2010**, *8(4)*, 1154–1177.
- 469 25. Sarich, M. Projected Transfer Operators. PhD thesis, Freie Universitt Berlin, 2011.
- 470 26. Sarich, M.; Schütte, C. Approximating Selected Non-dominant Timescales by Markov State
471 Models. *Comm. Math. Sci.* **2012**, *10*.
- 472 27. Schütte, C.; Noé, F.; Lu, J.; Sarich, M.; Vanden-Eijnden, E. Markov State Models Based on
473 Milestoning. *J. Chem. Phys* **2011**, *134 (19)*.
- 474 28. Faradjian, A.K.; Elber, R. Computing time scales from reaction coordinates by milestoning. *J.*
475 *Chem. Phys.* **2004**, *120*, 10880–10889.
- 476 29. Prinz, J.H.; Wu, H.; Sarich, M.; Keller, B.; Fischbach, M.; Held, M.; Chodera, J.D.; Schtte, C.;
477 Noé, F. Markov models of molecular kinetics: Generation and Validation. *J. Chem. Phys.* **2011**,
478 *134*, 174105.
- 479 30. Roeblitz, S. Statistical Error Estimation and Grid-free Hierarchical Refinement in Conformation
480 Dynamics. PhD thesis, FU Berlin, 2008.
- 481 31. Djurdjevac, N.; Sarich, M.; Schütte, C. On Markov state models for metastable processes.
482 *Proceeding of the ICM 2010 as invited lecture* **2010**. Download via
483 <http://www.math.fu-berlin.de/groups/biocomputing/publications/index.html>.
- 484 32. Horenko, I.; Dittmer, E.; Lankas, F.; Maddocks, J.; Metzner, P.; Schütte, C. Macroscopic
485 Dynamics of Complex Metastable Systems: Theory, Algorithms, and Application to B-DNA.
486 *J. Appl. Dyn. Syst.* **2009**.
- 487 33. Fleming, W.; Soner, H. *Controlled Markov Processes and Viscosity Solutions*; Springer; 2nd
488 edition, 2005.
- 489 34. Oksendal, B. *Stochastic Differential Equations*; Springer, 2003.
- 490 35. Pra, P.; Meneghini, L.; Runggaldier, W. Connections between stochastic control and dynamic
491 games. *Mathematics of Control, Signals and Systems* **1996**, *9*, 303–326.
- 492 36. Kushner, H.; Dupuis, P. *Numerical Methods for Stochastic Control Problems in Continuous Time*;
493 Springer Verlag, 1992.
- 494 37. Braack, M. *Finite Elemente*; 2012.
- 495 38. Banisch, R.; Hartmann, C. A meshfree discretization for optimal control problems. *to appear*
496 **2013**.
- 497 39. Sheu, S. Stochastic Control and Exit Probabilities of Jump Processes. *SIAM Journal on Control*
498 *and Optimization* **1985**, *23*, 306–328. [<http://epubs.siam.org/doi/pdf/10.1137/0323022>]
- 499 40. Latorre, J.; Metzner, P.; Hartmann, C.; Schtte, C. A Structure-preserving numerical discretization
500 of reversible diffusions. *Commun. Math. Sci.* **2011**, *9*, 1051–1072.

- 501 41. Chodera, J.D.; Elms, P.J.; Swope, W.C.; Prinz, J.H.; Marqusee, S.; Bustamante, C.; Noé, F.;
502 Pande, V.S. A robust approach to estimating rates from time-correlation functions. *arXiv:*
503 *1108.2304* **2011**. [[arXiv:cond-mat.stat-mech/1108.2304](https://arxiv.org/abs/cond-mat.stat-mech/1108.2304)]
- 504 42. Vanden-Eijnden, E. Transition Path Theory. *Lect. Notes Phys.* **2006**, *703*, 439–478.
- 505 43. Schtte, C.; Noe, F.; Lu, J.; Sarich, M.; Vanden-Eijnden, E. Markov State Models Based on
506 Milestoning. *J. Chem. Phys.* **2011**, *134* (20), 204105.

507 © June 2, 2013 by the authors; submitted to *Entropy* for possible open access
508 publication under the terms and conditions of the Creative Commons Attribution license
509 <http://creativecommons.org/licenses/by/3.0/>.