

# Collective Langevin dynamics of conformational motions in proteins

Oliver F. Lange and Helmut Grubmüller<sup>a)</sup>*Department of Theoretical and Computational Biophysics, Max-Planck-Institute for Biophysical Chemistry, Am Fassberg 11, 37077 Göttingen, Germany*

(Received 30 September 2005; accepted 3 April 2006; published online 2 June 2006)

Functionally relevant slow conformational motions of proteins are, at present, in most cases inaccessible to molecular dynamics (MD) simulations. The main reason is that the major part of the computational effort is spent for the accurate description of a huge number of high frequency motions of the protein and the surrounding solvent. The accumulated influence of these fluctuations is crucial for a correct treatment of the conformational dynamics; however, their details can be considered irrelevant for most purposes. To accurately describe long time protein dynamics we here propose a reduced dimension approach, collective Langevin dynamics (CLD), which evolves the dynamics of the system within a small subspace of relevant collective degrees of freedom. The dynamics within the low-dimensional conformational subspace is evolved via a generalized Langevin equation which accounts for memory effects via memory kernels also extracted from short explicit MD simulations. To determine the memory kernel with differing levels of regularization, we propose and evaluate two methods. As a first test, CLD is applied to describe the conformational motion of the peptide neurotensin. A drastic dimension reduction is achieved by considering one single curved conformational coordinate. CLD yielded accurate thermodynamical and dynamical behaviors. In particular, the rate of transitions between two conformational states agreed well with a rate obtained from a 150 ns reference molecular dynamics simulation, despite the fact that the time scale of the transition ( $\sim 50$  ns) was much longer than the 1 ns molecular dynamics simulation from which the memory kernel was extracted. © 2006 American Institute of Physics.

[DOI: [10.1063/1.2199530](https://doi.org/10.1063/1.2199530)]

## I. INTRODUCTION

Conformational motions of proteins are fundamental to protein function.<sup>1</sup> In recent years molecular dynamics (MD) simulations became more and more capable of elucidating functional processes of biomolecular systems.<sup>2</sup> However, MD operates in the full  $3N$ -dimensional configurational space of the protein and the surrounding solvent molecules (where  $N$  is the number of atoms). Consequently, the large number of pairwise interactions to be evaluated and the short time steps enforced by the fastest motions entail very long computation times in order to sufficiently sample the relatively slow conformational motions, which limits MD simulations at present to systems of  $10^5$ – $10^6$  atoms and to time scales of several 100 ns. Unfortunately, apart from a few exceptions, relevant biological processes, such as the gating of ion channels, allosteric interactions, ligand binding, molecular recognition, chemomechanical energy conversion, and many more, occur on microsecond to second time scales and therefore are currently far out of reach for conventional MD.

This holds true despite considerable efforts to speed up the computations, particularly of the long-range Coulomb forces, which have resulted in very efficient methods such as multiple step algorithms,<sup>3–8</sup> fast multipole methods,<sup>9–12</sup> and Ewald summation techniques.<sup>13</sup> Also, the use of

constraints<sup>14–16</sup> helps us to increase efficiency. Still, however, processes on time scales of microseconds and beyond can only be studied by resorting to certain “tricks” to enhance sampling of the conformational motions, as reviewed in Ref. 17. Unfortunately accelerated sampling necessarily implies loss of dynamical information and often, as for the faster methods, loss of thermodynamical accuracy as well.

To advance the methodology beyond conventional, “brute force” MD, a drastic *reduction of the large number of degrees of freedom* is therefore called for. This implies two steps. First, to identify few appropriate slow degrees of freedom<sup>18</sup> which serve to define a reduced active space within which the dynamics is evolved, without explicit treatment of the remaining orthogonal fast degrees of freedom. Second, to derive suitable equations of motion for these slow degrees of freedom.

Phenomenologically motivated selections of the active space include implicit solvent,<sup>19</sup> combined atom or bead models,<sup>20–22</sup> and the treatment of polypeptides as chains of stiff “platelets,” for which only  $\psi$ - $\phi$  backbone angles are retained as explicit degrees of freedom.<sup>23,24</sup> A somewhat related approach is the Gaussian network model.<sup>25</sup> However, by restriction to relevant atoms or groups of atoms and omitting others, only a very small subset of all possible collective degrees of freedom is considered. One may, therefore, expect to derive improved dimension-reduced descriptions of protein dynamics by dropping this empirical restriction and considering as degrees of freedom  $m$  fully general functions

<sup>a)</sup>Author to whom correspondence should be addressed. Fax: +49-551-201-2302. Electronic mail: [hgrubmu@gwdg.de](mailto:hgrubmu@gwdg.de)

$$c_i = f_i(\mathbf{x}_1, \dots, \mathbf{x}_N), \quad i = 1 \dots m, \quad (1)$$

of the atomic positions  $\mathbf{x}_j$ . Linear  $c_i$

$$c_i = \sum_{j=1}^N \mathbf{a}_{ij}(\mathbf{x}_j - \mathbf{x}_j^0) = \mathbf{A}(\mathbf{x} - \mathbf{x}_0) \quad (2)$$

are widely considered, e.g., within the framework of principal component analysis, which is often used to systematically derive slow and relevant (essential) collective degrees of freedom from MD simulations or structural ensembles.<sup>26–28</sup> Here we consider both linear and nonlinear collective degrees of freedom.

The question of how to reduce the full dynamics of all atomic degrees of freedom to equations of motion which describe the dynamics of the selected (usually slow) collective degrees of freedom has been addressed in the projection operator formalism by Zwanzig and Mori within a quite different context,<sup>29,30</sup> which led to the generalized Langevin equation (GLE).<sup>31–36</sup> Combining these two concepts, generalized collective degrees of freedom and dimension-reduced dynamics, we here develop the framework of collective Langevin dynamics (CLD) which describes protein dynamics in collective coordinates. The projection operator formalism is used to derive the necessary parameters for the GLE, i.e., an appropriate potential of mean force and memory kernels from short MD simulations. Thereby, all parameters are systematically obtained from first principles, and are specific for the chosen molecular system and the selected set of collective coordinates, which allows to automate parameter extraction.

The main tasks which need to be addressed are (1) identification of suitable conformational coordinates, (2) extraction of memory kernels from MD simulations, (3) construction of a suitable free energy landscape, and (4) evaluation of CLD accuracy and performance.

*Extraction of slow conformational degrees of freedom* with principal component analysis is well established. However, its separation of time scales has not yet been systematically assessed. For example, for the protein crambin it has been shown that principal component analysis (PCA) yields good separation of time scales (if and only) if all nonhydrogen atoms are included within the covariance matrix [OL, HG, submitted]. Whether principal components extracted from short MD simulations are able to describe protein dynamics on long time scales sufficiently well is, however, not clear, and subject to ongoing discussions.<sup>37,38</sup> We therefore revisited this question in more detail and found substantial evidence that indeed low-dimensional PCA subspaces extracted from short (~5 ns) MD simulations of T4 lysozyme and crambin describe more than 90% of the total atomic displacements observed in extended (>200 ns) MD simulations. Further evidence for the stability of PCA subspaces on long timescales is gained by showing that a 30 dimensional PCA subspace extracted for T4 lysozyme also allows an accurate description ( $\approx 1$  Å differences) of an ensemble of 38 crystallographic structures of T4 lysozyme [OL, HG, submitted].

For the first application of CLD reported here, we restricted the treatment to a single nonlinear conformational

coordinate, which is constructed in Sec. IV A based on principal component analysis and refined by human intervention.

In contrast, *extraction of memory kernels* from MD simulations is still a challenging problem. Despite considerable efforts,<sup>39–42</sup> a generally accepted approach has not yet emerged. Thus, we have studied and applied two memory extraction schemes and evaluated their performance within the framework of CLD.

The *free energy surface* can be determined by molecular dynamics sampling. More efficiently, however, are enhanced sampling techniques, for instance, multicanonical methods [e.g., replica exchange MD (REMD)],<sup>43</sup> smart Monte Carlo (SMC),<sup>44</sup> or umbrella sampling.<sup>45,46</sup> These techniques are complementary to CLD, because they yield canonical ensembles, but do not yield dynamical information. CLD, on the other hand, yields proper dynamical information, but relies on already known canonical ensembles.

*Assessment of the quality* of the obtained dimension-reduced description turns out to be nontrivial either. Clearly, direct comparison of the observed CLD trajectory  $\mathbf{c}(t) = c_1(t) \dots c_m(t)$  with explicit (deterministic) MD simulations is not meaningful, because the underlying GLE describes a stochastic process and because the dynamics is chaotic. Rather, suitable observables such as averages over many realizations of the stochastic process, or time averages such as time correlation functions, transport coefficients, or transition rates should be used.<sup>47</sup>

As we will demonstrate, time correlation functions are indeed well reproduced by CLD. However, in our view, they do not represent a rigorous test of CLD, because time correlation functions are closely related to the memory kernels that are used as input. Here, we use conformational transition rates as observables instead, which are fully unrelated to the input—yet statistically meaningful.

The Arrhenius equation<sup>48</sup>

$$k = \eta e^{-\beta \Delta G^\ddagger},$$

where  $\beta$  denotes the inverse temperature, clarifies in which way transition rates are influenced by the CLD parameters. The importance of the height of the free energy barrier  $\Delta G^\ddagger$  is evident immediately. Important too, however, is the prefactor  $\eta$ , which accounts for attempt frequency, recrossing events, and nonequilibrium effects. The height of the free energy depends on the choice of the conformational coordinate only, whereas the prefactor depends on the correct description of memory effects by the CLD model. Therefore, the check of the transition rates was also used to evaluate the relative performance of the different approaches to extract memory kernels.

For this comparison, conformational transition rates have to be calculated from a long reference MD simulation. Transitions which are suitable for this test have to be sufficiently fast, such that a number of transition events can be counted and their rate obtained with reasonable statistical accuracy. This is usually the case in small systems. Whereas the CLD framework derived below covers systems of arbitrary size, this consideration led us to choose the hexapeptide neurotensin as a first test system, which undergoes several folding/unfolding transitions at a 50 ns time scale.<sup>49</sup> Consid-

ering a curved conformational coordinate to describe these transitions we achieved a reduction of the high dimensional test system to a single degree of freedom.

The manuscript is organized as follows. In Sec. II A we derive the equations of motion for the general case of  $m$  nonlinear collective degrees of freedom. Section II B addresses the generalized friction terms of the GLE. The presented *memory equation* allows extraction of memory kernels from MD simulation, which determines the frictional force, and defines via the fluctuation dissipation theorem the statistical properties of the associated colored noise. In Sec. II C we apply Kramers's rate theory to CLD, as a reference for the numerical results.

After description of the computational methods, in Sec. IV CLD is applied to the hexapeptide neurotensin. Firstly, in Secs. IV A–IV C we will discuss the relevant conformational dynamics of neurotensin and its description with one curved conformational coordinate. Secondly, in Sec. IV D the CLD framework is used to extract model parameters from MD simulations. Finally, in Secs. IV E–IV G we discuss the performance of the one-dimensional CLD model of neurotensin focusing in Sec. IV F on a comparison of transition rates with explicit MD simulations of the full dynamics as a reference.

## II. THEORY

### A. Equations of motion for conformational dynamics

#### 1. Projection operator formalism

Let us first consider the conceptual framework, which we sketch here following<sup>50</sup> to clarify notation. We start with the full molecular dynamics, which are described by the Hamiltonian

$$\mathcal{H}(\mathbf{x}, \mathbf{p}) = \frac{1}{2} \sum_{i=1}^n \frac{p_i^2}{m_i} + \mathcal{V}(\mathbf{x}), \quad (3)$$

where  $\mathbf{x}$  and  $\mathbf{p}$ , with components  $x_i$  and  $p_i$ , respectively, are the  $n$ -dimensional position and momentum vectors and  $m_i$  their masses. A solution of the corresponding canonical equations is defined through an initial value  $(\mathbf{x}_0, \mathbf{p}_0)$ ; to each initial condition corresponds a trajectory,  $\varphi(t) = \varphi(\mathbf{x}_0, \mathbf{p}_0, t)$ . Subsequently the subscript 0 is omitted.

In the framework of the projection operator formalism a dynamical variable,<sup>29,30,51</sup> mechanical property,<sup>50</sup> or physical quantity<sup>52</sup> is defined as a mapping on phase space  $\mathbb{R}^{3N} \times \mathbb{R}^{3N}$

$$A: \mathbb{R}^{3N} \times \mathbb{R}^{3N} \rightarrow \mathbb{R}^{2m}$$

$$(\mathbf{q}, \mathbf{p}) \mapsto A(\mathbf{q}, \mathbf{p}),$$

with the  $2m$  components denoted by  $A_i$ ,  $i=1, \dots, 2m$ . The space  $\mathcal{D}$  of all dynamical variables is endowed with the inner products

$$\langle A|B \rangle_{ij} = \int A_i(\mathbf{x}, \mathbf{p}) B_j(\mathbf{x}, \mathbf{p}) \rho(\mathbf{x}, \mathbf{p}) d\mathbf{x} d\mathbf{p}.$$

Here  $\rho$  is the canonical distribution  $\rho(\mathbf{x}, \mathbf{p}) = Z^{-1} e^{-\beta \mathcal{H}(\mathbf{x}, \mathbf{p})}$  with partition function  $Z$  and inverse temperature  $\beta$ . We use the

bracket formalism and denote the elements of  $\mathcal{D}$  as *ket*-vectors  $|A\rangle$ .

A dynamical variable varies in time through its argument; a dynamical variable whose value at  $t=0$  was  $A(\mathbf{x}, \mathbf{p})$  acquires at time  $t$  the value  $A(\varphi(\mathbf{x}, \mathbf{p}, t))$ . One can also take a ‘‘Heisenberg’’ or ‘‘Lagrangian’’ point of view and introduce a time-dependent dynamical variable  $e^{\mathcal{L}t}A$ , where  $\mathcal{L}$  denotes the Liouville operator defined by the Poisson bracket

$$\mathcal{L} = \{ \cdot, \mathcal{H} \} = \sum_{i=1}^n \left( \frac{\partial \mathcal{H}}{\partial p_i} \frac{\partial}{\partial x_i} - \frac{\partial \mathcal{H}}{\partial x_i} \frac{\partial}{\partial p_i} \right). \quad (4)$$

The propagator  $e^{\mathcal{L}t}$  allows us to define the time-dependent *ket*-vector

$$|A(t)\rangle \equiv |e^{\mathcal{L}t}A\rangle = A(\varphi(\mathbf{x}, \mathbf{p}, t)),$$

which obeys the Liouville equation

$$\frac{d}{dt} |A(t)\rangle = \mathcal{L} |A(t)\rangle. \quad (5)$$

The projection operator<sup>30</sup>

$$\mathcal{P} = \langle A|A\rangle^{-1} |A\rangle \langle A| \quad (6)$$

allows us to separate the time dependence of the dynamical variable into a part within the linear subspace  $U$  spanned by the *ket*-vectors  $|A_i\rangle$  and a part within the orthogonal subspace  $U^\perp$ ,<sup>53</sup>

$$\frac{d}{dt} |A(t)\rangle = e^{\mathcal{L}t} \mathcal{P} \mathcal{L} |A\rangle + \int_0^t d\tau e^{\mathcal{L}(t-\tau)} \mathcal{P} \mathcal{L} |F(t)\rangle + |F(t)\rangle, \quad (7)$$

with the random force  $|F(t)\rangle \equiv e^{(1-\mathcal{P})\mathcal{L}t} (1-\mathcal{P}) \mathcal{L} |A\rangle$ . The random force lies within the subspace  $U^\perp$ , i.e.,  $(1-\mathcal{P})|F(t)\rangle = |F(t)\rangle$ , therefore

$$\begin{aligned} \mathcal{P} \mathcal{L} |F(t)\rangle &= \mathcal{P} \mathcal{L} (1-\mathcal{P}) |F(t)\rangle \\ &= \langle \mathcal{L} (1-\mathcal{P}) F(t) | A \rangle \langle A | A \rangle^{-1} |A\rangle. \end{aligned}$$

Defining the *memory function*  $\Gamma(t) \equiv \langle \mathcal{L} (1-\mathcal{P}) F(t) | A \rangle \times \langle A | A \rangle^{-1}$ ,<sup>39</sup> Eq. (7) becomes the GLE,

$$\frac{d}{dt} |A(t)\rangle = \mathcal{P} \mathcal{L} |A(t)\rangle + \int_0^t d\tau \Gamma(\tau) |A(t-\tau)\rangle + F(t). \quad (8)$$

Thus the dynamics of  $|A\rangle$  are split into the dynamics within  $U$  and a correction term which describes the evolution of the system in  $U^\perp$ .  $F(t)$  is the *random force*<sup>30</sup> exerted by the uncoupled motion in  $U^\perp$ , i.e.,  $\langle F(t) | A(0) \rangle = 0$ , with its realization depending on the chosen initial conditions for the orthogonal part of the motion. The energy uptake due to the random force  $F(t)$  is counterbalanced by the generalized friction, as expressed formally by the generalized fluctuation dissipation theorem

$$\langle F(t) | F(t') \rangle = \langle A | A \rangle \Gamma(t-t'). \quad (9)$$

#### 2. Definition of motion along conformational coordinate(s) as the observable

Here we propose to apply the projection operator formalism rigorously to the dynamics of suitably chosen collective

degrees of freedom, Eq. (1), and derive equations of motion for them. To be specific, and for simplicity of notation, we consider the dynamics of *one* nonlinear collective variable  $c$  ( $m=1$ ), although the theory can be generalized to more dimensions in a straightforward manner. The dynamics of the collective degree of freedom  $c$  are best represented by motion along a suitably chosen one-dimensional submanifold  $\mathcal{M} \subset \mathbb{R}^{3N}$  of the configurational space parametrized by  $c$ . However, in practice, at first a submanifold  $\mathcal{M}$  will be chosen which is able to represent the motion of interest, which defines the collective degree of freedom as a projection to that submanifold.

To derive the equations of motions for the collective coordinate with the projection operator formalism, the problem is recast in terms of a dynamical variable  $A$  with two components. The first component,  $A_1$ , is given by the projection of vector  $\mathbf{x}$

$$A_1: \Gamma \rightarrow \mathbb{R}$$

$$(\mathbf{x}, \mathbf{p}) \mapsto c = f(\mathbf{x}),$$

and the second component by the orthogonal projection of the momentum  $\mathbf{p}$  onto the tangential space  $T_{f(\mathbf{x})}\mathcal{M}$  of the manifold to the point corresponding to parameter  $f(\mathbf{x})$

$$A_2: \Gamma \rightarrow T_{f(\mathbf{x})}\mathcal{M}$$

$$(\mathbf{x}, \mathbf{p}) \mapsto \dot{c} = \nabla_{\mathbf{x}} f \cdot \mathbf{M}^{-1} \mathbf{p},$$

where  $\mathbf{M}$  is a diagonal mass matrix. For a one-dimensional equation of motion, a suitably chosen reduced mass  $\mu$  is required, which is derived from Eq. (10) via the equipartition theorem,  $\langle \dot{c}^2 \rangle = (\beta \mu)^{-1}$ . The mean squared velocity

$$\langle \dot{c}^2 \rangle = \int (\nabla_{\mathbf{x}} f \cdot \mathbf{M}^{-1} \mathbf{p})^2 \rho(\mathbf{x}, \mathbf{p}) d\mathbf{x} d\mathbf{p}$$

consists of a sum of pure terms  $\sim p_i^2$  and mixed terms  $\sim p_i p_j$ . After integration over the momenta the mixed terms vanish, which allows, via  $\int p_i^2 d\mathbf{p} = \beta^{-1} m_i$ , us to define the reduced mass as

$$\mu = \left( \int \sum_{i=1}^n \left( \frac{\partial f}{\partial x^i} \right)^2 \frac{1}{m_i} \rho(\mathbf{x}) d\mathbf{x} \right)^{-1}.$$

### 3. Equations of motion for conformational coordinate(s)

The above definitions allow for the application of the projection-operator formalism in order to derive the equations of motion for the collective degree of freedom(s). To this aim, and exploiting the fact that the two components of the dynamical variable  $A$  are conjugated variables, the system of the two first order GLEs, Eq. (8), is cast into one second order GLE. This is possible because the first components of random force and memory function vanish, which can be seen from

$$\mathcal{L}|A_1\rangle = \sum_i \frac{\partial \mathcal{H}}{\partial p_i} \frac{\partial c}{\partial x_i} \in T_{f(\mathbf{x})}\mathcal{M},$$

such that the orthogonal part  $(1-\mathcal{P})\mathcal{L}|A_1\rangle$  vanishes. Hence, the fluctuation-dissipation theorem, Eq. (9), simplifies accordingly to

$$\langle F_2(t)|F_2(t')\rangle = \langle A|A\rangle \gamma(t-t'), \quad (11)$$

with  $\gamma(t) \equiv \Gamma_2(t)$ .

The conventional way to proceed from here is to apply the linear projector  $\mathcal{P}$ , Eq. (6), to the remaining component  $\mathcal{P}\mathcal{L}|A_2(t)\rangle$  of the first term in Eq. (8), thus obtaining an effective harmonic force  $\Omega$ .<sup>30,54</sup> However, here we avoid this harmonic approximation by adopting the nonlinear projection operator originally introduced by Zwanzig<sup>51</sup> and used recently by Chorin *et al.*<sup>53</sup> Apart from introducing a dependency of  $\gamma$  on the ket-vector  $|A\rangle$ , this generalization does not change the above derivation, Eqs. (7) and (8).<sup>53</sup>

To be able to project to a curved conformational coordinate we generalize the operator defined in Ref. 53 to

$$\mathcal{P}|\cdot\rangle = \frac{1}{\rho^c(c, \dot{c})} \int |\cdot\rangle \rho(\mathbf{x}, \mathbf{p}) d\Omega(c, \dot{c}), \quad (12)$$

with  $d\Omega(c, \dot{c}) := \delta_c \delta_{\dot{c}} (\nabla_{\mathbf{x}} f \cdot \mathbf{M}^{-1} \mathbf{p} - \dot{c}) d\mathbf{x} d\mathbf{p}$ . Here, we have defined the *conformational density*  $\rho^c(c)$  (Ref. 55) as the projection of the density in configurational space onto the conformational coordinates,

$$\rho^c(c, \dot{c}) = \int \rho(\mathbf{x}, \mathbf{p}) \|\nabla_{\mathbf{x}} f \mathbf{M}^{-1/2}\|^2 d\Omega(c, \dot{c}). \quad (13)$$

The mass-matrix  $\mathbf{M}$  would be rendered to unity, if mass-weighted coordinates ( $\tilde{\mathbf{x}} = \mathbf{M}^{1/2} \mathbf{x}$  and  $\tilde{\mathbf{p}} = \mathbf{M}^{-1/2} \mathbf{p}$ ) were used. Since the chosen example comprises masses in the range of 12–16 amu and, therefore, the difference is small, we have not used mass-weighted coordinates here.

Applying the generalized projector shows that  $\mathcal{P}\mathcal{L}|A_2(t)\rangle$  is the expectation value of the potential force acting tangentially to  $\mathcal{M}$  under all possible realizations of the trajectory and a correction term due to the curvature of the chosen parametrization  $f$  of  $\mathcal{M}$ ,

$$\begin{aligned} \mathcal{P}\mathcal{L}|A_2(t)\rangle &= \frac{1}{\rho^c(c, \dot{c})} \int [\nabla_{\mathbf{x}} \mathcal{V} \cdot \nabla_{\mathbf{x}} f \mathbf{M}^{-1} \\ &\quad - \nabla_{\mathbf{p}} \mathcal{H} \cdot \nabla_{\mathbf{x}} (\nabla_{\mathbf{x}} f \cdot \mathbf{M}^{-1} \mathbf{p})] \rho(\mathbf{x}, \mathbf{p}) d\Omega(c, \dot{c}). \end{aligned} \quad (14)$$

Defining the potential of mean force as  $W(c, \dot{c}) = -\beta^{-1} \log \rho^c(c, \dot{c})$ , we show in the Appendix that its derivative

$$\frac{\partial W}{\partial c}(c, \dot{c}) = -\beta^{-1} \frac{1}{\rho^c(c, \dot{c})} \frac{\partial \rho^c(c, \dot{c})}{\partial c} \quad (15)$$

yields the right hand side of Eq. (14). In the linear case, integration on the momentum part of the integral can be carried out separately, such that the dependence on the velocity  $\dot{c}$  vanishes. This yields the final result

$$\mathcal{P}\mathcal{L}|A_2(t)\rangle = \frac{1}{\mu} \frac{\partial W}{\partial c}(c). \quad (16)$$

For reasons of practicality we approximate in the nonlinear case by averaging out the dependence on the velocities.

To cast the resulting equation into the more usual form<sup>31,56,57</sup> of a second order GLE, we set  $R(t) = \mu F_2(t)$ , and from Eq. (8) one obtains

$$\mu \ddot{c}(t) = -\frac{d}{dc}W(c(t)) - \int_0^t d\tau \mu \gamma(t-\tau, c) \dot{c}(\tau) + R(t), \quad (17)$$

which is, except for the approximation in Eq. (16) for nonlinear  $f$ , the *exact* equation of motion for the projected dynamics. Its right hand side is composed of a potential of mean force  $W$ , a generalized friction  $\gamma$ , and a random force  $R$ . The latter two obey the corresponding fluctuation-dissipation theorem,

$$\langle R(0)|R(t)\rangle = \mu \beta^{-1} \gamma(t).$$

The computation of the random force  $R(t)$  requires solution of a Liouville equation which is far more complicated than the original unprojected problem. The advantage of the reformulation of the equations of motion in the form of the GLE Eq. (17) is, of course, that the random force can be replaced by a stochastic term, i.e., a randomly generated force with similar statistical properties. In particular, its autocorrelation function has to satisfy the fluctuation-dissipation theorem Eq. (9). Accordingly, the remainder of this section addresses the task to extract the three components  $W$ ,  $\gamma$ , and  $R$  from atomistic molecular dynamical simulations.

## B. Extraction of Langevin parameters from MD simulations

The potential of mean force  $W(c) = -\beta^{-1} \log \rho^c(c)$  is obtained from the configurational density projected to the chosen collective coordinate(s). Here, the necessary canonical ensembles will be generated by MD simulations, although any method that yields a canonical ensemble can be used, e.g., (REMD),<sup>43</sup> umbrella sampling,<sup>45,46</sup> or metadynamics.<sup>58</sup>

As described within Sec. II B 1, the generalized friction is obtained via velocity autocorrelation functions, which in turn are obtained from MD simulations. Here relatively short trajectories contain already sufficient information, because the respective memory kernels typically decay rapidly. The memory kernels obtained in that way determine the statistical properties of the random force via the fluctuation-dissipation theorem. Therefore, no additional parameters are needed, accordingly we address within the second subsection further below, the generation of instances of a random force process  $R$  from a given autocorrelation function.

### 1. Memory equation

The memory equation<sup>39,59,60</sup>

$$\frac{d}{dt}\Psi(t) = - \int_0^t d\tau \gamma(t-\tau)\Psi(\tau) \quad (18)$$

connects the velocity autocorrelation function (VACF),  $\Psi(t) = \langle \dot{c}(0)|\dot{c}(t)\rangle / \langle \dot{c}^2 \rangle$ , with the memory kernel,  $\gamma(t)$ . In its usual form Eq. (18) is obtained from the GLE Eq. (17) setting the potential  $W \equiv 0$  by application of  $\mu^{-1} \langle \dot{c}^2 \rangle^{-1} \langle \dot{c}(0)|$  and noting that  $\langle F(t)|A(0)\rangle = 0$ .

$\Psi(t)$  can be computed readily from MD simulations, such that Eq. (18) has to be solved to obtain the memory kernel  $\gamma(t)$ . Note, however, that this usual form of the memory equation,<sup>39,59,60</sup> when applied to dynamics under the influence of deterministic forces  $W(c)$ , yields an adulterated  $\gamma(t)$ . These additional forces are misattributed to the total random and frictional force, which should exclusively determine  $\gamma(t)$ . We therefore suggest to consider this contribution separately. The additional term takes the form of a velocity/force correlation,

$$\Pi(t) \equiv \mu^{-1} \langle \dot{c}^2 \rangle^{-1} \left\langle \dot{c}(0) \frac{\partial W}{\partial c}(c(t)) \right\rangle, \quad (19)$$

and serves to quantify the accuracy of the usual approximation Eq. (18).

Indeed, in the following application to neurotensin  $\Pi(t)$  was found to be some magnitudes smaller than the other terms involved such that Eq. (18) is used throughout the following discussion, but an extension of the presented methods to incorporate  $\Pi(t)$ , if becoming necessary, should be straightforward.

Note that the treatment with the usual form of the memory equation is consistent, if the memory is interpreted or used together with a GLE where  $W \equiv 0$ . However, even in these cases it is impractical and counterintuitive, because the influence of the unaccounted forces can alter a memory kernel such that it does not decay to zero anymore. For example, a motion in a harmonic potential  $W = 0.5\omega c^2$  can accurately be described by a GLE with  $W \equiv 0$ , but asymptotically  $\gamma(\infty) = \omega/m$ .

For completeness, we note that an alternative memory equation can be obtained following Berkowitz *et al.*<sup>61</sup> by applying  $\mu^{-1} \langle \dot{c}(0)|$  to the GLE, Eq. (17), which leads in our case to

$$\begin{aligned} \langle \dot{c}(0)|\dot{c}(t)\rangle &= \mu^{-1} \left\langle c(0) \frac{\partial W}{\partial c}(c(t)) \right\rangle - \int_0^t d\tau \gamma(t-\tau) \\ &\quad \times \langle c(0)|\dot{c}(t)\rangle. \end{aligned} \quad (20)$$

However, we do not consider this any further, because it contains slowly converging positional contributions to the autocorrelation functions.

### 2. Random force

To generate instances of  $R(t)$  from  $\gamma(\tau)$  via the fluctuation-dissipation theorem, we follow the method proposed in Ref. 56, which is exact, in contrast to other methods.<sup>57,62,63</sup> Briefly, the Wiener-Khinchin theorem is exploited, which connects the spectral density

$$J(\omega) = \int_{-\infty}^{\infty} dt \langle R(0) | R(t) \rangle e^{-i\omega t} \quad (21)$$

to the power spectrum by

$$J(\omega) = \left| \int_{-\infty}^{\infty} dt R(t) e^{-i\omega t} \right|^2. \quad (22)$$

Hence, the average amplitude of the Fourier transformed noise  $\langle |R(\omega)| \rangle$  is determined by the memory function  $\gamma(t) = \mu\beta^{-1} \langle R(0) | R(t) \rangle$ . An otherwise random process, whose Fourier amplitudes are, however, fixed in that way, is gained by setting

$$R(t) = \int_{-\infty}^{\infty} d\omega \sqrt{J_K(\omega)} z(\omega) e^{i\omega t},$$

where  $z(\omega) \in \mathbb{C}$  are realizations of a normal distribution with unit variance, and  $J_K$  is the spectral density corresponding to the given memory function

$$J_K(\omega) = m\beta \int_{-\infty}^{\infty} dt \gamma(t) e^{-i\omega t}.$$

This method allows for the generation of noise for any given autocorrelation function.

### C. Transition rates with Kramers's theory

To check how accurately the dimension-reduced GLE approximates the fully atomistic dynamics, we will compare the conformational transition rates of the CLD model with the rates obtained from explicit MD simulations. The transition rates for the conformational dynamics governed by the GLE of the CLD model can be obtained in two ways. Either the GLE is integrated numerically, which yields a trajectory whose transitions can be counted, or transition rates are estimated directly from the GLE using Kramers's theory.<sup>48</sup> For the latter we follow Kramers's approach and approximate the potential of mean force with parabolas at the minima  $W_\alpha(x) \approx \mu\omega_\alpha^2 x^2$  and at the barrier top  $W_\ddagger(x) \approx -\mu\omega_\ddagger^2 x^2$ . Then the escape rate is

$$k_\alpha = \left( \sqrt{\frac{\hat{\gamma}^2(\xi)}{4} + \omega_\ddagger^2} - \frac{\hat{\gamma}(\xi)}{2} \right) \frac{\omega_\alpha}{2\pi\omega_\ddagger} \exp(-\beta\Delta W_\alpha^\ddagger), \quad (23)$$

with index  $i=A,B$  for states  $A$  and  $B$ , respectively, and  $\Delta W_\alpha^\ddagger = W^\ddagger - W_\alpha$  the height of the barrier. Here  $\hat{\gamma}(z)$  denotes the Laplace transform of the memory kernel  $\gamma(t)$ , and  $\xi$  is subject to the condition

$$\xi = -\frac{\gamma(\hat{\xi})}{2} + \sqrt{\frac{\gamma(\hat{\xi})}{4} + \omega_\ddagger^2}. \quad (24)$$

In the case of memory-free friction,  $\gamma(t) = 2\gamma_{\text{eff}}\delta(t)$ , Eq. (23) simplifies due to  $\hat{\gamma}(\xi) \equiv \gamma_{\text{eff}}$ , and adopts the widely known form.<sup>48</sup> For a comprehensive review, we refer to Ref. 48.

## III. METHODS

### A. Molecular dynamics simulation

All molecular dynamics (MD) simulations were carried out using the GROMACS simulation suite.<sup>64</sup> Lincs and Settle<sup>15,16</sup> were applied to constrain covalent bond lengths, allowing an integration step of 2 fs. Electrostatic interactions were calculated using the particle-mesh-Ewald method.<sup>65,66</sup> The temperature was kept constant by separately coupling ( $\tau=0.1$  ps) the peptide and solvent to an external temperature bath.<sup>67</sup> The pressure was kept constant by weak isotropic coupling ( $\tau=0.1$  ps) to a pressure bath.<sup>67</sup>

Several molecular dynamics simulations of neurotensin were carried out, using the optimized potentials for liquid simulations all atom force field.<sup>68</sup> Neurotensin, a peptide with the sequence Ac-RRPYIL, was solvated with 2246 TIP4P water molecules and 2  $\text{Cl}^-$  counterions in a cubic box. A first simulation was started from an extended configuration and equilibrated for 10 ns. A 90 ns simulation (NT1) was started from the last snapshot of the equilibration, and coordinates were recorded every 1 ps. A second simulation (NT2) with a length of 63 ns was started from the last snapshot of NT1, and positions and velocities were recorded every 10 fs, which allowed for the computation of velocity autocorrelations without aliasing artifacts.

Additionally, eight 500 ps simulations,  $\text{NTS}_i$ ,  $i=1\cdots 8$ , were started from snapshots of NT1 selected for mutually large root mean square differences, and positions and velocities were recorded every 10 fs.

### B. Principal component analysis

Principal component analyses (PCAs) were carried out by diagonalizing the covariance matrix  $\mathbf{C} = \langle (\mathbf{x} - \langle \mathbf{x} \rangle) (\mathbf{x} - \langle \mathbf{x} \rangle)^T \rangle$ , where  $\mathbf{x}$  denotes protein or peptide atomic positions in the  $3N$ -dimensional configurational space. Translational and rotational motions were removed by least squares fitting to a reference structure  $\mathbf{x}_{\text{ref}}$ ,<sup>26</sup> the angular brackets denote the averages over a MD trajectory. The eigenvectors of  $\mathbf{C}$  yielded the PCA modes  $\{\mathbf{a}_j\}_{j=1\cdots 3N}$ , and positions projected onto mode  $j$  were obtained as  $c_j = \mathbf{a}_j \cdot (\mathbf{x} - \langle \mathbf{x} \rangle)$ .

For consistency with the positions, the projected velocities  $\dot{c}_j(t) = \mathbf{a}_j \cdot \mathbf{v}_c(t)$  were computed from corrected velocities  $\mathbf{v}_c(t)$  without contribution of translational and rotational motions such that  $c_j(t) = \int_0^t \dot{c}_j(\tau) d\tau + c_j(0)$ . For this correction translational and rotational velocities were computed using the displacement vector  $d(t_i)$  that moves the center of mass into the origin, and the rotation matrix  $R(t_i)$  which minimized root mean square deviation (RMSD) to the reference structure, obtained from the least square fitting of the positions. Thus, the corrected velocities were given by

$$\mathbf{v}_c(t_i) = \mathbf{v}(t_i) - \Delta t [d(t_{i-1}) - d(t_i) + R(t_{i-1})\mathbf{x}(t_i) - R(t_i)\mathbf{x}(t_i)],$$

where  $\Delta t$  denotes the sampling interval.

### C. Definition of a one-dimensional curved conformational coordinate

By visual inspection of the projection of trajectory NT1 onto the first three eigenvectors of  $\mathbf{C}$ , eight snapshots

$\{\mathbf{x}_{\text{sel},i}\}_{i=1\dots 8}$  were selected evenly spaced along the observed trace of high conformational density (Fig. 2). To remove bias introduced by choosing single snapshots out of a large number of equally reasonable alternatives, averages over all  $n_i$  snapshots  $\{\mathbf{x}_{\text{sel},i}^j\}_{j=1\dots n_i}$  within a sphere of radius 0.1 nm [in the three-dimensional (3D) projection] around  $\mathbf{x}_{\text{sel},i}$  were used. The conformational coordinate was then constructed by cubic spline interpolation between these averages in the full  $3N$ -dimensional space. Subsequent discretization yielded  $N = 1200$  points  $\{\mathbf{z}_i\}_{i=1\dots N}$ .

## D. Projection onto the conformational coordinate

The conformational coordinate defined by the discretized submanifold  $\mathcal{M} = \{\mathbf{z}_i\}_{i=1\dots N}$  was parametrized by a mapping function  $f$  [see Eq. (1)], such that  $c = f(\mathbf{z}_1) = 0$  and  $c = f(\mathbf{z}_{1200}) = 1$ , respectively. All intermediate values were defined via the contour length  $s_j = \sum_{k=2}^j \|\mathbf{z}_k - \mathbf{z}_{k-1}\|$  as  $c = f(\mathbf{z}_j) := s_j / s_N$ . Thus, the length unit,  $L$ , of the projected coordinate is  $L = s_N$ , and the metric of the configurational space is preserved upon projection.

Unfortunately, the straightforward approach to project a point  $\mathbf{x}$  onto the point of  $\mathcal{M}$  which is closest in space

$$P(\mathbf{x}) = \arg \min_{z \in \mathcal{M}} \|\mathbf{x} - z\| \quad (25)$$

led to several “wrong” projections due to the U shape of the coordinate. In particular, and as will be discussed in Sec. IV, this simple projection scheme, therefore, resulted in unphysical discontinuities.

This problem was resolved by additionally considering the time information of the trajectory. Specifically, snapshots close in time were enforced to yield projections close to each other. To determine the projection  $P(\varphi(\mathbf{x}, \mathbf{p}, t_i))$ , we proceeded as follows. First, both the discretized conformational coordinate  $\mathbf{z}_i$  and the trajectory  $\varphi(\mathbf{x}, \mathbf{p}, t_i)$  were projected preliminary onto the first 100 principal modes (obtained as above) yielding  $\bar{\mathbf{z}}_i$  and  $\bar{\varphi}(\mathbf{x}, \mathbf{p}, t_i)$ , respectively. Second, the final projection of the trajectory to the curved coordinate was determined via

$$P(\varphi(\mathbf{x}, \mathbf{p}, t_{i+1})) = \arg \min_{z \in I(c(t_i), r)} \|\bar{\varphi}(\mathbf{x}, \mathbf{p}, t_{i+1}) - \bar{\mathbf{z}}\|, \quad (26)$$

where the interval of the conformational coordinate

$$I\{c(t_i), r\} = \{z \in \mathcal{M} | c(t_i) - r \leq f(z) \leq c(t_i) + r\}$$

defines a window of width  $2r$  centered around the previous result of the projection  $c(t_i) = f[P(\varphi(\mathbf{x}, \mathbf{p}, t_i))]$ . Parameter values below 0 or greater than 1 were allowed by extending the conformational coordinate linearly at both ends. The window size  $r$  was chosen to trade off sufficient fast response of the projection with robustness against unphysical jumps; for the 10 fs sampling,  $r = 1/1200$  and for the 1 ps sampling  $r = 1/12$ . Velocities were projected onto the first 100 principal modes as described above, and subsequently onto the tangent to the conformational coordinate at the point  $P(\mathbf{x}(t_i))$ .

## E. Potential of mean force

To compute the potential of mean force  $W$  along the conformational coordinate,  $W(c) = -kT \log \rho(c)$ , the density

$\rho(c)$  was obtained from the MD ensemble projected to the conformational coordinate. For this purpose, histograms with 100 bins were determined and smoothed by convolution with a Gaussian function of width  $\sigma = 0.025L$ , where  $L$  denotes the length scale of the conformational coordinate. Outside the sampled range of  $c$  the potential was continued by a harmonic potential  $W_{\text{harm}}(c) = 11.5(c - 0.5L)^2$  as

$$W_{\text{extend}}(c) = [1 - S(c)]W(c) + S(c)W_{\text{harm}}(c),$$

where the switching function is defined by the sigmoidal function  $S(c) = \{1 + \exp[-50(c - 1)]\}^{-1} + 1 - \{1 + \exp[-50c]\}^{-1}$ . Forces were computed by linear interpolation between the numerically obtained derivatives at neighboring discretization points.

## F. Solution of the memory equation

Memory kernels  $\gamma(t)$  were obtained by solving the Volterra equation of the first kind, Eq. (18),

$$\frac{d}{dt}\Psi(t) = - \int_0^t d\tau \gamma(t - \tau)\Psi(\tau), \quad (27)$$

with the velocity autocorrelation function  $\Psi(t) = \langle \dot{c}(0) | \dot{c}(t) \rangle / \langle \dot{c}^2 \rangle$  computed from the MD simulation.

Integration of this equation is notoriously unstable, because the result  $\gamma$  does not depend continuously on  $\Psi$ .<sup>69</sup> For the case at hand  $\Psi$  contains statistical noise, which aggravates the problem. Therefore, the often used approach to transform the equation into a Volterra equation of the second kind<sup>69,70</sup> does not yield the hoped for continuous dependence on  $\Psi$ , because it appears in both terms, the convolution term and on the left hand side. Furthermore, the differentiation usually increases the noise level considerably, which also implies instabilities.<sup>70</sup>

To find physically meaningful solutions we therefore resorted to regularizations. This can be done either by imposing a model function for the result or by imposing local criteria, e.g., smoothness. Here we tested these two main approaches. Very strong regularization was achieved with the first approach by imposing a model function with only three free parameters and weak regularization was used for the second approach, where sequential Tikhonov regularization<sup>71</sup> was applied to favor smooth solutions.

For the first approach, which would subsequently be denoted by FIT, we used the three-parameter model function

$$\gamma_{\text{fit}} = 2\gamma^f \delta(t) + Ae^{-at}. \quad (28)$$

With this ansatz, the memory equation can be solved analytically by using Laplace transformations. The Laplace transform of the memory equation

$$z\hat{\Psi}(z) - 1 = -\hat{\gamma}(z)\hat{\Psi}(z),$$

together with that of the memory kernel  $\hat{\gamma}(z) = \gamma^f + A/(z + a)$ , yields the transformed velocity autocorrelation function  $\hat{\Psi}(z) = (z + \gamma^f + A/(z + a))^{-1}$ , whose back transformation is

$$\Psi_{\text{fit}}(t) = \exp(-(a + \gamma^f)t/2) \times \left[ \cosh(Rt/2) + \frac{a - \gamma^f}{R} \sinh(Rt/2) \right], \quad (29)$$

with  $R := \sqrt{a^2 - 2a\gamma^f + (\gamma^f)^2 - 4A}$ . To obtain the parameters  $\gamma^f$ ,  $a$ , and  $A$  of Eq. (28), Eq. (29) was least squares fitted to the numerically obtained velocity autocorrelation functions,  $\Psi(t_i)$ , with the curve fitting tool of MATLAB(tm). Note that for a negative radicant an imaginary  $R$  results, which might be difficult to handle by common fitting algorithms. Nevertheless, a purely real form of Eq. (29) can be obtained for these cases by replacing  $R$  with  $i\bar{R}$ , which renders the cosh and sinh into cos and sin, respectively.

For the second approach, which will subsequently be denoted by DIR, we discretized ( $\Delta t = t_i - t_{i-1} = 0.01$  ps) the memory equation accounting explicitly for a delta-function-like contribution to the memory [cf. Eq. (28)]

$$\dot{\Psi}_i = -\gamma^f \Psi_i - \Delta t \sum_{k=0}^{i-1} \gamma_k \Psi_{i-k} \omega_{ik} - \gamma_i \Psi_0 \omega_{ii} \Delta t. \quad (30)$$

This equation was solved using an adaption of sequential Tikhonov regularization,<sup>71</sup> as described in Ref. 72, which finds the regularized solution as minimum of the Tikhonov criterion

$$\min \left\{ \sum_{i=0}^N (I_i(\gamma^f, \gamma_0, \gamma_1, \dots) - \dot{\Psi}_i)^2 + \alpha^2 \Omega(\gamma^f, \gamma_0, \gamma_1, \dots)^2 \right\}, \quad (31)$$

where  $I_i(\gamma^f, \gamma_0, \gamma_1, \dots)$  denotes the right hand side of Eq. (30), and  $\Omega$  denotes the approximation to the second derivative of  $\gamma$ . A large regularization parameter  $\alpha$  favors smooth solutions at the cost of a larger residual norm, while a small  $\alpha$  has the opposite effect.

The regularization parameter can be chosen from an analysis of the  $L$  curve.<sup>73-75</sup> However, here the  $L$ -curve optimum of roughly  $\alpha=20$  was not very pronounced. In order to contrast the strong regularization method above with a method, which does not bias the result too much, we chose with the help of the  $L$  curve the relatively low regularization parameter  $\alpha=0.14$ .

For illustration purposes we also obtained memory kernels  $\gamma_{1\text{-reg}}$  and  $\gamma_{11\text{-reg}}$  with  $\alpha=20$ , but these memory kernels were not used for the CLD model.

## G. Langevin simulation

The GLE was integrated via the adaption of the velocity Verlet scheme<sup>76</sup> described in the work of Tuckerman and Berne,<sup>77</sup> in particular, Eqs. (3.1)–(3.7) therein. The most time-consuming part of the computation is the calculation of the generalized friction term of Eq. (8), because it involves summation over many previous velocities, i.e.,

$$(\gamma * \dot{c})_i = \sum_{k=0}^i \omega_k \gamma_k \dot{c}_{i-k}, \quad (32)$$

with  $i=0 \cdots N$ .

Because the velocity  $\dot{c}_i$  is computed iteratively, the usual advantage of the fast Fourier transform (FFT) method to compute the convolution  $(\gamma * \dot{c})_i$  for all values of  $i$  simultaneously<sup>78,79</sup> cannot be exploited readily. Nevertheless, a strategy by Baker and Derakhshan<sup>80</sup> was applied, which allows for FFT treatment of large parts of the convolution sum. In this way only a small part involving the newest velocities has to be summed explicitly every time step.<sup>85</sup>

For the integration of the corresponding (nongeneralized) Langevin equation we set  $\omega_0 \gamma_0 = 2\gamma^f / \Delta t$  and  $\gamma_k = 0$  for  $k > 0$ , and replaced the random forces by  $R_n = (2kTm\gamma^f / \Delta t)^{-1/2} \xi_n$ , where the  $\xi_n$  are independent Gaussian random variables with zero mean and  $\langle \xi_n^2 \rangle = 1$ , and  $\Delta t$  denotes the integration time step.

## H. Statistics of conformational transitions

Transition rates were determined from the one-dimensional projection of molecular dynamics trajectories to the conformational coordinate or from the one-dimensional CLD trajectories by counting. First, every snapshot  $c(t_i)$  was assigned to one of the two conformational states  $s(t_i) = A, B$  and then the number of changes of  $s$  was evaluated. To account for nonthermalized recrossings<sup>48</sup> a variable threshold was applied to the low-pass filtered projection  $\tilde{c}(t_i)$ , which depended on the previous conformation  $s(t_{i-1})$

$$s(t_i) = \begin{cases} A & s(t_{i-1}) = A \wedge \tilde{c}(t_i) < 0.66 \\ B & s(t_{i-1}) = B \wedge \tilde{c}(t_i) > 0.36. \end{cases}$$

As low-pass filter we used smoothing with a Gaussian function of width  $\sigma=40$  ps. The transition rate  $k_{\alpha\beta}$  for the transition  $\alpha \rightarrow \beta$  was given by  $k_{\alpha\beta} = n_{\alpha\beta} / (N_\alpha \Delta t)$ , where  $n_{\alpha\beta}$  denotes the changes of  $s(t_i)$  from state  $\alpha$  to state  $\beta$  and  $N_\alpha$  is the number of snapshots for which  $s(t_i) = \alpha$ .

The threshold value and the bandwidth of the low-pass filter were chosen manually and introduce clearly a bias into the obtained rates. However, here we only need to compare the rates obtained with the same method, such that this bias was canceled out. Moreover, other sets of parameters tested did not change the relative differences between CLD and reference transition rates.

Confidence intervals were determined via the Poisson-statistic  $P_\lambda(n) = e^{-\lambda} \lambda^n / n!$ , since transitions were rare events. Via  $\langle n^2 \rangle = \lambda = \langle n \rangle$  the number of observed transitions  $n$  determined the Poisson-parameter  $\lambda$  and with that an estimate of the error of the transition rate. A 95% confidence interval in the logarithmic representation was computed by choosing its width  $d$ , such that  $P_\lambda(k \in [n \exp(-d), n \exp(d)]) = 0.95$ . In the case of a large number of observed transitions  $n > 60$  the Poisson statistics was approximated by the error function via  $P_\lambda(a \leq k \leq b) = \Phi(b/\sqrt{\lambda}) - \Phi(a/\sqrt{\lambda})$ , with  $\Phi(x) = \int_{-\infty}^x (2\pi)^{-1/2} \exp(-x^2/2) = \frac{1}{2} [\text{erf}(x/\sqrt{2}) + 1]$ .

## IV. RESULTS AND DISCUSSION

The framework of collective Langevin dynamics (CLD) has been laid out in the Theory section, where the equations of motion for slow collective coordinates were derived.



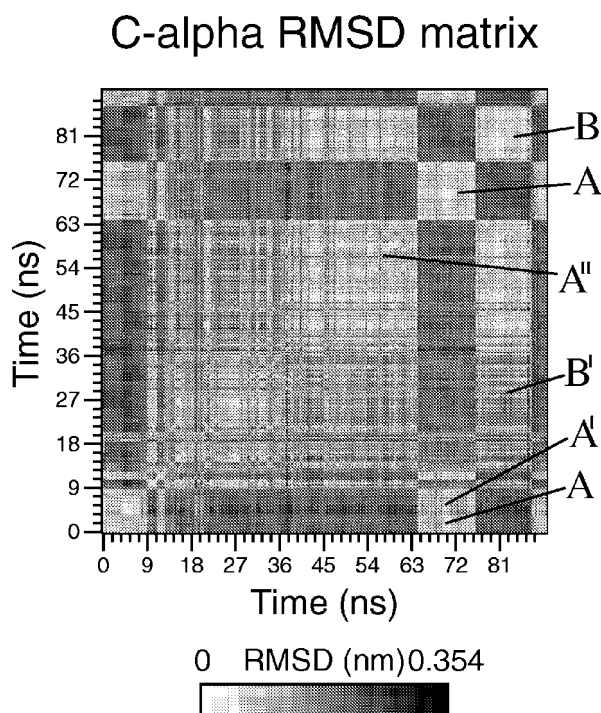


FIG. 1. Root mean square deviation (RMSD) of the  $C_{\alpha}$  atoms for each pair of snapshots of trajectory NT1. The RMSD ranges from 0 (white) to 0.382 nm (dark). The labels indicate conformational substates, see text.

In this section we will show that CLD describes the conformational dynamics of a model protein sufficiently well to allow reliable prediction of transition rates. To this end, we will compare the transition rates obtained from CLD models with reference transition rates obtained by a standard MD simulation. Here, a relatively small test system had to be used, because for the accurate rate estimates the transitions need to occur sufficiently often within a MD simulation. We therefore choose the hexapeptide neurotensin, which, due to fast conformational dynamics, underwent several transitions already in a 150 ns MD simulation.

As a first step we modeled the CLD of neurotensin by means of a one-dimensional coordinate. Whereas from the methodological point of view this appears to be the simplest case, reduction from  $3N$  coordinates to a single one is of course the most drastical case possible and, hence, represents a hard test. To this end we had to use a curved coordinate, as will be described in Sec. IV A. Subsequently, free energy and generalized friction will be extracted for the chosen coordinate from explicit MD simulation, and, finally, transition rates will be compared between MD simulations and the CLD model.

### A. Neurotensin as a test system

Having described all parts of the CLD framework, we now will apply it by considering the conformational dynamics of neurotensin as a specific test case. This system has been chosen, because we expected it to undergo sufficiently many conformational transitions at the MD time scale to allow comparisons of transition rates.

Indeed, as can be seen in Fig. 1, neurotensin underwent several main conformational transitions  $A \rightarrow B$  during the

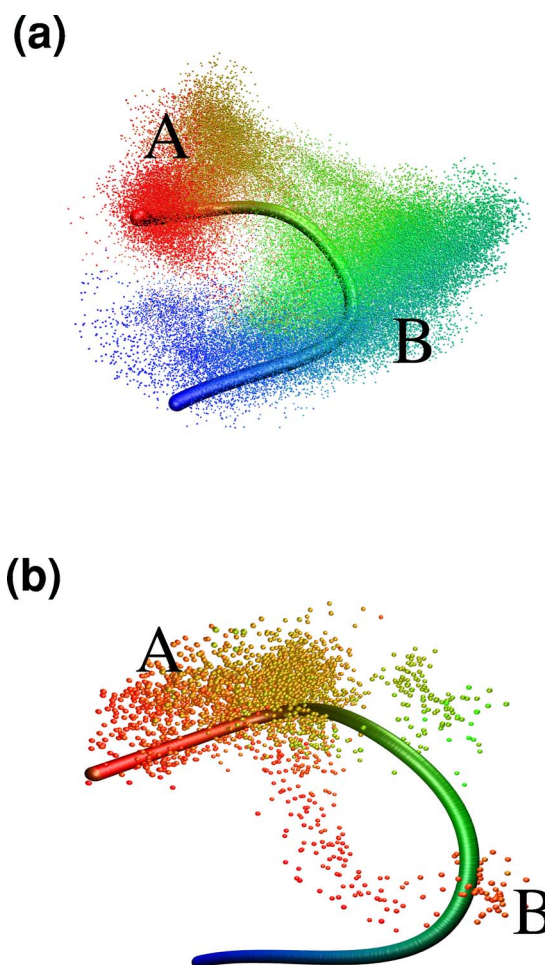


FIG. 2. Projection of the conformational coordinate (thick line) and the configurational ensemble (dots) onto the first three PCA modes. The colors denote the resulting mapping of snapshots  $\mathbf{x}_i$  to position on the coordinate  $c=f(\mathbf{x}_i)$ , from red,  $c \approx 0$ , to blue,  $c \approx 1$ . (a) Whole configurational ensemble NT1. (b) Interval (70–75 ns) of NT1, where a substate of A is visited.

90 ns MD simulation, NT1. The figure shows the matrix of the root mean square deviation (RMSD) of the  $C_{\alpha}$  atoms for each pair of snapshots of the trajectory NT1. Conformational states were defined as almost invariant subsets of the configurational space.<sup>81</sup> They are visible in the RMSD matrix as distinct bright blocks on the diagonal. Bright off-diagonal blocks indicate that a certain conformational substate was revisited. Interestingly, as can also be seen in the figure, the two main conformational states subdivide further into substates (denoted by primes) as is typical for proteins,<sup>82</sup> thus giving rise to a complex conformational dynamics also within the main states. In this sense, the system represents a particularly harsh test system for CLD: The CLD model has to predict correct first passage times without knowledge of the substate dynamics. This lack of knowledge is of course intrinsic to a dimension-reduced approach and it is important to find out how well CLD can cope with it.

### B. Construction of a curved conformational coordinate

As a first task, we need to construct a collective coordinate, which resolves both states A and B. We start by analyzing the MD ensemble, as projected onto the first three

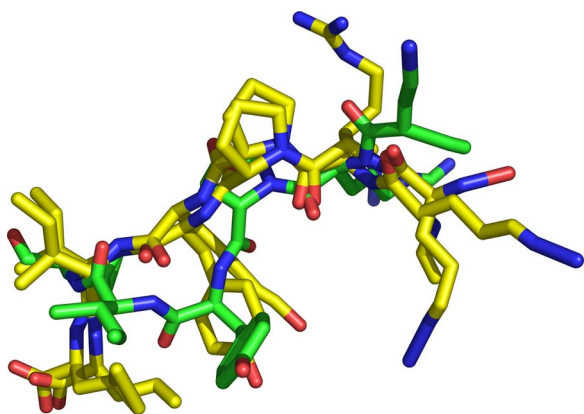


FIG. 3. Overlay of average neurotensin structures. The relative orientation of the structures minimizes the RMSD between the  $C_\alpha$  atoms. The green structure is obtained from state A, and the two yellow structures are obtained from state B. The parts of the side chains that were overly distorted due to the averaging were removed. The N terminus is oriented towards the upper right corner.

principal components, shown in Fig. 2(a). The red points represent the structures belonging to conformation A, the blue points belong to conformation B, and the green points belong to transitions between both states. Some structures representative for these conformational states are shown in Fig. 3. The shape of this ensemble was such that no conceivable

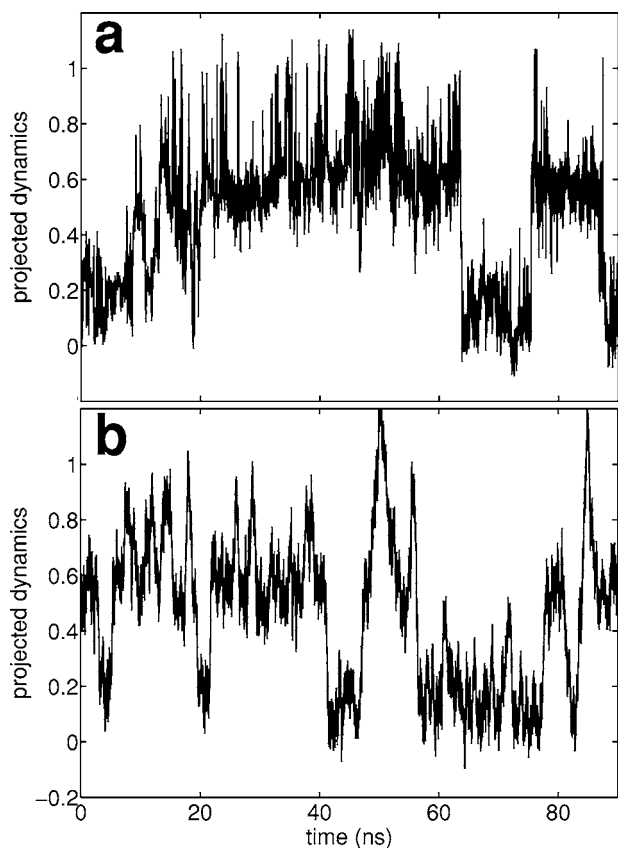


FIG. 4. (a) Projection of neurotensin internal motion onto conformational coordinate. This shows the projection of NT1 onto the conformational coordinate. A number of transitions occur between the clearly distinguishable two states centered at 0.2 and 0.6, respectively. (b) An example of a CLD trajectory. The plot shows a 90 ns of a trajectory generated by the model  $CLD_{1-fit}$  (cf. Sec. IV E).

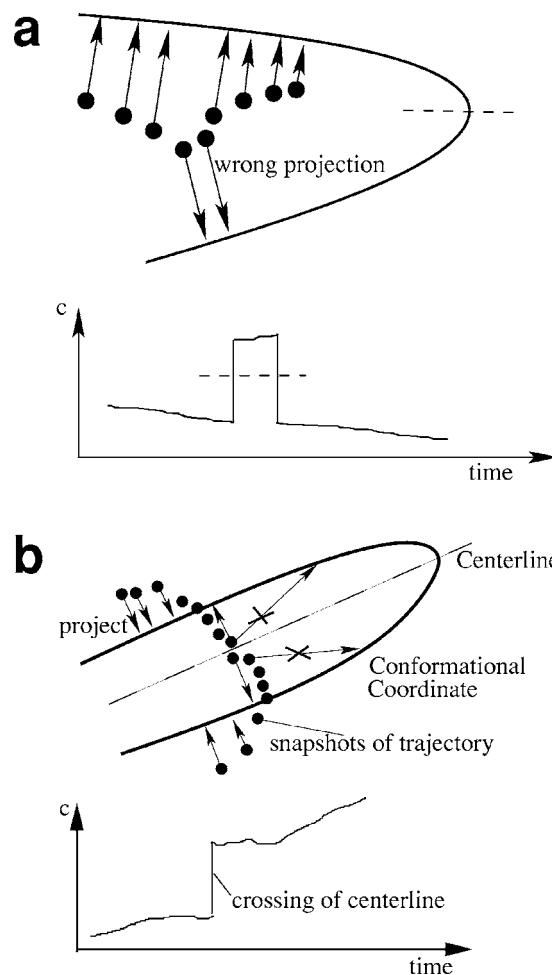


FIG. 5. Illustration of main sources for artifacts in the projection to a curved coordinate. Snapshots in the configurational space (circles) are projected (arrows) to a curved coordinate by a pure distance criterion. The resulting projection is plotted against time under the pictures. (a) A trajectory moves from right to left along one arm of the coordinate, i.e., the projection is decreasing (cf. plot). However, two snapshots are slightly closer to the other side of the curved coordinate and, hence, the projection erroneously jumps to large values and back, although no real conformational transition occurred. (b) In contrast to (a), here a real transition from the low projection part to the high projection part of the coordinate occurs. However, the trajectory crosses to far away from the curved coordinate and, therefore, shortcuts the bulge drastically, which results in an artifactual large jump in the projection, in the moment of crossing of the centerline.

able *linear* coordinate would resolve the two conformational states. In particular, the close blue and red points are separated by a free energy barrier, which cannot be resolved by a linear coordinate. We therefore constructed the curved coordinate shown in Fig. 2(a) (see Sec. III), which clearly resolves states A and B. The projection of trajectory NT1 onto the coordinate  $c$  (cf. Fig. 4) revealed several well resolved transitions between the conformational substates A and B, centered around  $c \approx 0.2$  and  $c \approx 0.6$ , respectively. This projection turned out not to be straightforward, and care had to be taken to avoid possible artifacts.

The more technical aspects of this projection described below are not of direct relevance for the CLD model; we have included a brief description, nonetheless, to illustrate problems, which typically arise from the use of curved coordinates as well as their solutions.

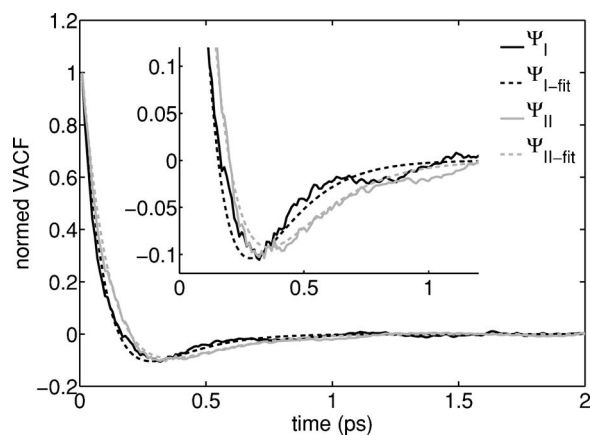


FIG. 6. Velocity autocorrelation function (VACF) of MD trajectories (solid) and their respective fits to Eq. (29) (dashed). The inset shows the same data enlarged.

The main problem arose from the fact that no *direct* transitions between the two main states, *A* and *B*, were seen in the vicinity of the red and blue points [cf. Fig. 2(a)], but only *indirect* ones in the region of the green points. Therefore, straightforward assignment of each MD structure to the nearest point of the conformational coordinate would fabricate transitions, as illustrated in Fig. 5(a). These spurious transitions would adulterate the reference transition rates used for confirmation of the CLD model. This problem was solved by taking time information into account (cf. Sec. III). Careful inspection showed that the spurious transitions were indeed eliminated.

Figure 2(b) shows an extreme example. Here, several structures seem to approach state *B* in the projection onto the first three principal components. Accordingly, these structures would be assigned to state *B* in any purely distance based projection onto the shown curved coordinate. However, as can also be seen in Fig. 1, the RMSD to state *A* remains small for all shown structures and large to state *B*, such that assignment to *B* would be wrong. Indeed, as seen from the coloring in Fig. 2(b), all snapshots of the substate of *A* were correctly allocated to conformation *A*.

Figure 5(b) illustrates and explains a second problem (see caption), resulting in discontinuities in the projected motion. This problem was solved by careful placement of the curved coordinate.

### C. Velocity autocorrelation function of collective motion

The velocity autocorrelation function (VACF) is required to derive the memory kernels for the CLD model and, therefore, needs to be extracted from the MD trajectories. Further below we will analyze how well this observable is reproduced by the CLD model.

Two VACFs,  $\Psi_I$  and  $\Psi_{II}$ , were obtained from trajectories NT2 and NT2<sub>II</sub>, respectively (see Sec. III). NT2<sub>II</sub> refers to the interval 10–19 ns of the 63 ns trajectory NT2. Both VACFs are shown in Fig. 6, together with their respective fits to Eq. (29),  $\Psi_{I\text{-fit}}$  and  $\Psi_{II\text{-fit}}$ . Both  $\Psi_I$  and  $\Psi_{II}$  are well ap-

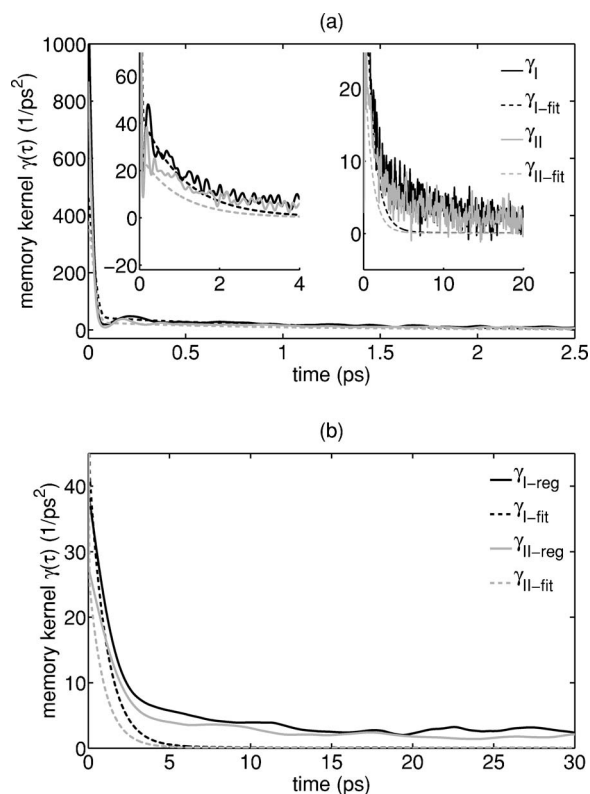


FIG. 7. Memory kernel functions computed from the VACFs shown in Fig. 6. (a) Memory computed with DIR (solid) and with FIT (dashed). The insets show the same data in different zooms. (b) Memory kernels from the same VACFs, but method DIR was used with a higher regularization parameter to eradicate oscillations.

proximated by the fits and are very similar with respect to the position and depth of the first minimum, which indicates good convergence of their main feature.

Their rapid decay shows that most correlations occur at a picosecond time scale. Furthermore, the pronounced negative dip in intermediate time scales ( $0.3 \text{ ps} < \tau < 0.5 \text{ ps}$ ) indicates resonant behavior or memory effects in the system. The similarity to the dip in VACFs of simple liquids caused by caging of the tagged molecule by its immediate neighbors<sup>59</sup> is suggestive.

The medium scale oscillations of the VACFs, which are seen to modulate the dominating features captured by the fit (cf. inset of Fig. 6), however, indicate more complex dynamics than typically observed for simple liquids. For example, the slowly decaying oscillatory contributions to the VACF are clearly above the noise threshold seen for larger times  $\tau > 5 \text{ ps}$ . Although these medium scale modulations are due to relatively fast dynamics (i.e., on the picosecond time scale), the difference between  $\Psi_I$  and  $\Psi_{II}$  indicates that this feature may not be fully converged, which suggests a weak dependence of the VACF on the much slower conformational degree(s) of freedom. Note that also correlations on much slower time scales exist (see Sec. IV G) but are dominated by correlations on fast time scales, such that for  $\tau > 5 \text{ ps}$  they are not discernable from noise.

### D. Extraction of memory kernels

From the VACF, we can now proceed and compute the memory kernel as the essential quantity that captures the

TABLE I. Parameters of velocity autocorrelation function, Eq. (29), obtained for different trajectories. The presented parameters for NTS<sub>3</sub> and NTS<sub>4</sub> were the most extreme of all 500 ps trajectories. The parameters for unions, e.g., NTS<sub>1</sub>+NTS<sub>5</sub>, were obtained by fitting to an overlay of the VACFs of the respective trajectories. The labels NTS<sub>uneven</sub>, NTS<sub>even</sub>, and NTS<sub>all</sub> denote NTS<sub>1</sub>+NTS<sub>3</sub>+NTS<sub>5</sub>+NTS<sub>7</sub>, NTS<sub>2</sub>+NTS<sub>4</sub>+NTS<sub>6</sub>+NTS<sub>8</sub>, and NTS<sub>1</sub>+...+NTS<sub>8</sub>, respectively.

|                                    | $a$ (ps <sup>-1</sup> ) | $\gamma$ (ps <sup>-1</sup> ) | $A$ (ps <sup>-2</sup> ) | mass $m$ (amu) |
|------------------------------------|-------------------------|------------------------------|-------------------------|----------------|
| NT2                                | 0.78                    | 14.8                         | 49.2                    | 7.3            |
| NT2 <sub>II</sub>                  | 1.1                     | 11.5                         | 27.0                    | 8.9            |
| NTS <sub>4</sub>                   | 1.06                    | 11.7                         | 28.1                    | 13.33          |
| NTS <sub>3</sub>                   | 1.45                    | 12.5                         | 43.9                    | 12.94          |
| NTS <sub>1</sub> +NTS <sub>5</sub> | 1.04                    | 12.6                         | 33.6                    | 11.1           |
| NTS <sub>2</sub> +NTS <sub>6</sub> | 1.32                    | 13.2                         | 35.5                    | 9.84           |
| NTS <sub>3</sub> +NTS <sub>7</sub> | 1.24                    | 13.1                         | 35.0                    | 12.52          |
| NTS <sub>4</sub> +NTS <sub>8</sub> | 1.07                    | 12.2                         | 30.7                    | 8.78           |
| NTS <sub>uneven</sub>              | 1.13                    | 12.9                         | 34.3                    | 11.82          |
| NTS <sub>even</sub>                | 1.18                    | 12.7                         | 32.9                    | 9.31           |
| NTS <sub>all</sub>                 | 1.16                    | 12.8                         | 33.6                    | 10.56          |

influence of the many degrees of freedom excluded from explicit treatment in the CLD model. To this aim the memory equation was here solved using two different methods, FIT and DIR (cf. Sec. III F).

Figure 7(a) compares the memory kernels  $\gamma_I$  and  $\gamma_{II}$  computed with DIR with the respective memory kernels computed with FIT,  $\gamma_{I\text{-fit}}$  and  $\gamma_{II\text{-fit}}$ . As described in Sec. III, method FIT admits only a certain type of functions for the memory kernel, and hence involves stronger regularization constraints than DIR, which allows any sufficiently smooth function.

Overall, all memory functions are very similar. In particular, they drop rapidly to ca. 5% of their initial values at  $\tau \approx 0$ , followed by a decay with a 1 ps time constant  $a$  (cf. Table I). Closer inspection reveals also small deviations [cf. enlargement in the right inset of Fig. 7(a)]. In particular, the memory kernels  $\gamma_I$  and  $\gamma_{II}$ —obtained with the less regularizing method DIR—show strong oscillations, a second slower decay component, and do not approach zero. None of these features is seen in the memory kernels  $\gamma_{I\text{-fit}}$  and  $\gamma_{II\text{-fit}}$ . These features, therefore, deserve a closer analysis.

As can be seen from the left inset in Fig. 7(a), most details of the fast oscillations vary for the different trajectories. They are due to the unconverged medium scale oscillations and the small scale statistical noise of the VACF, both

TABLE II. The first two columns show effective friction constants estimated from input VACFs ( $\Psi_I$ ,  $\Psi_{II}$ ,  $\Psi_{I\text{-fit}}$ , or  $\Psi_{II\text{-fit}}$ ) or from corresponding memory functions ( $\gamma_I$ ,  $\gamma_{II}$ ,  $\gamma_{I\text{-fit}}$ , or  $\gamma_{II\text{-fit}}$ ). The second two columns show forward and backward transition rates observed in trajectories of the respective CLD models. The reference transition rates from the MD trajectory are provided in the first line.

|                       | $\int \gamma dt$ (ps <sup>-1</sup> ) | $1/\int \Psi$ (ps <sup>-1</sup> ) | $k_{A \rightarrow B}$ (10 <sup>-4</sup> ps <sup>-1</sup> ) | $k_{B \rightarrow A}$ (10 <sup>-4</sup> ps <sup>-1</sup> ) |
|-----------------------|--------------------------------------|-----------------------------------|------------------------------------------------------------|------------------------------------------------------------|
| MD                    |                                      |                                   | 1.9 +1.5/-0.9                                              | 1.4 +1.1/-0.6                                              |
| CLD <sub>I</sub>      | 249                                  | <147                              | 0.3 +0.7/-0.2                                              | 0.2 +0.5/-0.2                                              |
| CLD <sub>II</sub>     | 187                                  | <120                              | 0.7 +0.7/-0.4                                              | 0.5 +0.5/-0.2                                              |
| CLD <sub>I-fit</sub>  | 49                                   | 56                                | 2.2 +1.1/-0.7                                              | 1.4 +0.7/-0.5                                              |
| CLD <sub>II-fit</sub> | 42                                   | 30                                | 2.1 +1.2/-0.8                                              | 1.0 +0.6/-0.4                                              |

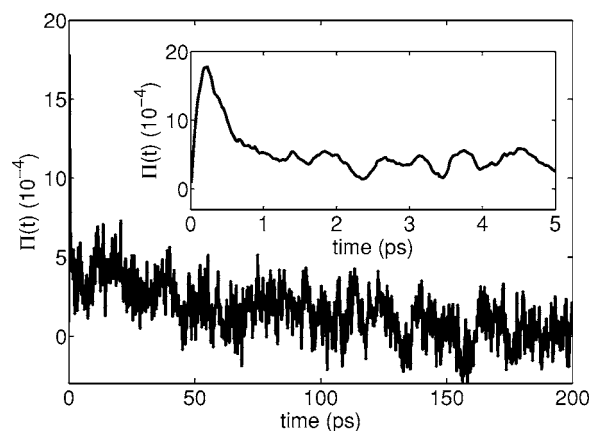


FIG. 8. Potential term  $\Pi(t)$  of memory equation. The inset shows the same data enlarged.

strongly amplified by the inherent instability of the memory equation. Accordingly, the large amplitude of the oscillations in the memory kernels likely has no physical basis. To aid the remaining discussion, Fig. 7(b) also shows memory kernels, whose oscillations were removed by increasing the regularization parameter  $\alpha$  as defined in Eq. (31).

In contrast to these oscillations, both remaining features not seen in  $\gamma_{I\text{-fit}}$  and  $\gamma_{II\text{-fit}}$ , the slower decay component and the lack of complete decay to zero for very long times, are comparable for both memory functions  $\gamma_{I\text{-reg}}$  and  $\gamma_{II\text{-reg}}$ , as shown in Fig. 7(b). Contrary to what one might expect on first sight, this “amplified noise” does not reflect long time memory effects. For  $\tau \gg 5$  ps,  $\gamma(\tau)$  is not well defined by the memory equation, because  $\Psi(\tau)$  is dominated by noise for these longer times. Therefore, setting  $\gamma(\tau) \equiv 0$  for these long times satisfies the memory equation equally well and removes this purely numerical artifact.

That long time memory effects are overestimated by DIR is further supported by the relation for the effective friction constants,<sup>59</sup>  $(\int_0^\infty \Psi(\tau) d\tau)^{-1} = \int_0^\infty \gamma(\tau) d\tau$ . Indeed, the effective friction constants estimated from the shown interval of  $\gamma_I$  and  $\gamma_{II}$ , respectively, are significantly higher than those derived from the corresponding VACFs (cf. Table II), suggesting a spuriously slow decay of the memory kernels. We also note that the numerical evaluation of  $(\int_0^\infty \Psi(\tau) d\tau)^{-1}$  does not converge and, hence, the values reported in the table give only an upper bound. Therefore, the real discrepancy between the effective frictions is even larger.

To quantify the effect of the artificial long tails of the memory kernels on the observed dynamics, we obtained a new set of memory kernels  $\gamma_{I\text{-tail}}$  and  $\gamma_{II\text{-tail}}$  by manually damping the tail to zero beyond 20 ps. As will be discussed in Sec. IV F these new memory kernels yield improved accuracy for the transition rates. The particular choice of the cut-off time  $\tau_c$  does not influence the results significantly; e.g., similar results were obtained for  $\tau_c = 10$  ps. Note that this procedure removes only the incomplete decay, but leaves the slow decay component of the memory kernel between  $0 < \tau < 10$  ps unchanged. We, therefore, attribute this remaining difference to the memory kernels obtained with FIT to a physical basis.

We finally note that in the present context the term  $\Pi(\tau)$ ,

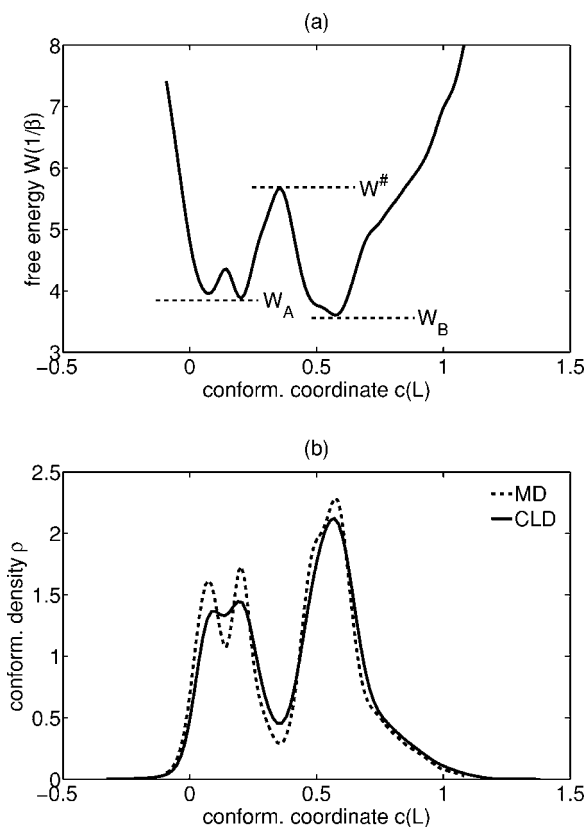


FIG. 9. (a) Potential of mean force along the conformational coordinate. The energy levels depicted as  $W_A$ ,  $W_B$ , and  $W^\#$  were used for calculation of barrier heights in Kramers's theory. (b) Comparison of conformational density of a CLD ensemble with that of the reference MD ensemble.

as defined by Eq. (19), of the memory equation was neglected. This term, derived from a correlation function between mean force and velocity, corrects for those velocity correlations, which are caused by the inertial motion of the system within a nonzero free energy surface rather than by memory effects due to the eliminated degrees of freedom. Due to the highly diffusive nature of the conformational dynamics of the system at hand, here the influence of the free energy on the velocities is small and, therefore, also the term  $\Pi(\tau)$  is expected to be small, implying that it can be neglected to good approximation. Indeed, as shown in Fig. 8,  $\Pi(\tau)$  is three orders of magnitude smaller than the VACF term for small  $\tau$ , and for larger times  $\tau > 5$  ps it is one magnitude smaller than the noise in the VACF, which justifies our approximation.

### E. Conformational dynamics by CLD

In the following three sections we test how well the dynamics along the conformational coordinate is actually described by the CLD model. Additionally to the memory kernels obtained above, a reduced mass and a free energy are required.

The free energy [Fig. 9(a)] along the conformational coordinate  $c$  was obtained from the conformational density [Fig. 9(b)] as potential of mean force, averaged over both available MD ensembles, NT1 and NT2. The reduced mass

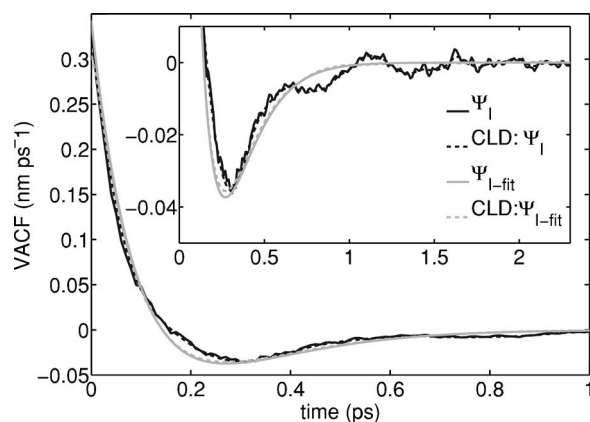


FIG. 10. Comparison of CLD generated (dashed lines) with reference MD (solid lines) velocity autocorrelation functions. Note that the dashed lines are hardly seen, since the respective curves match very good.

$\mu$  was obtained via the equipartition theorem  $\langle \dot{c}^2 \rangle = (\beta\mu)^{-1}$  from the amplitude of the velocity fluctuations (cf. Table I).

Having thus obtained all parameters directly from MD simulations, collective Langevin dynamics trajectories of 300 ns length each were obtained by numerical integration of the generalized Langevin equation, Eq. (17).

In the following we analyze the accuracy of the CLD model in terms of suitable dynamical and thermodynamical observables of the CLD model.

Firstly, we compare the thermodynamic properties to those obtained from the reference MD simulation. Since all thermodynamic observables of this CLD model can be obtained from its one-dimensional partition function, it suffices to compare the conformational density  $\rho$  with that of the MD ensembles, projected to the conformational coordinate [cf. Fig. 9(b)]. As can be seen, the densities agreed well with each other, although that of the CLD model was slightly smoother. This result confirms that the used friction kernel and random forces generated from it satisfy the fluctuation dissipation theorem.

Secondly, the dynamics was checked by comparison of the VACF with references from the MD simulations.

We focus on the evaluation of the CLD models based on  $\Psi_I$ , because the results for the CLD models obtained from  $\Psi_{II}$  were similar. Figure 10 shows the VACFs of the two CLD models using  $\gamma_I$  and  $\gamma_{I\text{-fit}}$  together with the reference VACF of the MD. All VACFs agree well. In particular, the initial decay and the position of the dip were well reproduced.

As was expected, the VACF obtained with  $\gamma_{I\text{-fit}}$  is nearly identical with the fit to the MD VACF  $\Psi_{I\text{-fit}}$ . Therefore, for this model the quality of the fit determines the accuracy. This restriction is gone if the method DIR is used to obtain the memory kernels. As shown in the inset of Fig. 10 the resulting VACF reproduces the reference very closely. This was quantified by the deviation  $(\int (\Psi - \Psi_{\text{CLD}})^2 dt)^{1/2}$  between CLD-VACF and reference, which was smaller for DIR ( $4.8 \times 10^{-3}$  nm/ps) than for FIT ( $8.2 \times 10^{-3}$  nm/ps).

### F. Prediction of transition rates by CLD

It was shown that CLD yields trajectories with accurate conformational densities and VACFs. Although these proper-

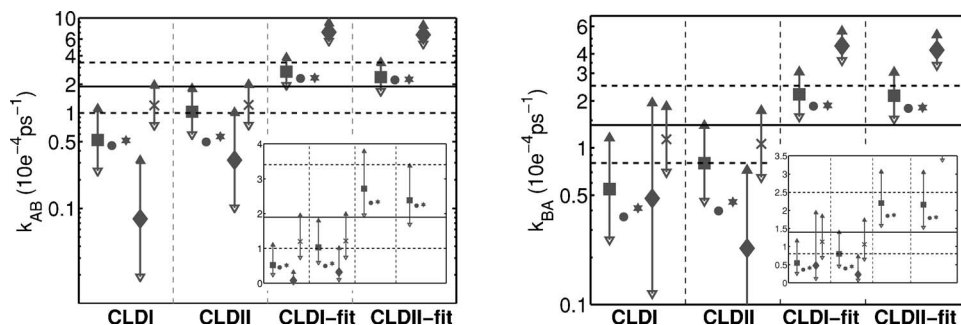


FIG. 11. Comparison of transition rates. The reference transition rate with its 95% confidence interval is shown by the solid and slashed horizontal lines. The CLD transition rates with error bars denoting 95% confidence intervals (if available) are grouped in the order  $CLD_I$ ,  $CLD_{II}$ ,  $CLD_{I-fit}$ , and  $CLD_{II-fit}$ . For every model four rates were obtained: stochastic simulation (squares), memory-free and full-memory Kramers's theory (circles and stars, respectively), and stochastic simulation of the memory-free Langevin equation (diamonds). Additionally, rates for  $CLD_I$  and  $CLD_{II}$  obtained by stochastic simulation with the shortened memory functions  $\gamma_{I-tail}$  and  $\gamma_{II-tail}$ , respectively, are denoted by crosses.

ties provide a useful consistency check, we do not consider them as a rigorous test of CLD, because they were also used as input for the CLD model. In contrast, transition rates were not used for the parametrization. As the most rigorous test, we therefore finally check forward and backward transition rates against references obtained from a long MD simulation.

In the following, we label the results from the different approaches as  $CLD_I$ ,  $CLD_{II}$ ,  $CLD_{I-fit}$ , and  $CLD_{II-fit}$ . The first two denote the CLD model whose memory was obtained with DIR from  $\Psi_I$  and  $\Psi_{II}$ , respectively, and the latter two denote the corresponding CLD models whose memory was obtained with FIT. In Fig. 11 transition rates observed from 300 ns CLD trajectories are shown (squares) with error bars indicating their 95% confidence interval (cf. Sec. III). For comparison, the reference transition rates obtained from MD simulations NT1 and NT2 with a total simulation time of 153 ns are shown as horizontal lines.

Additionally, transition rates corresponding to the respective CLD models were estimated from Kramers's theory. This estimate relies on the generalized Langevin equation of CLD in harmonic approximation to the free energy. The respective curvatures were determined at the minima as  $c_A = 65(\beta L^2)^{-1}$ ,  $c_B = 76(\beta L^2)^{-1}$  and at the barrier as  $c_{\ddagger} = 210(\beta L^2)^{-1}$  by fitting parabolas to the free energy profile (cf. Fig. 9). The barrier heights were  $W_{\ddagger} - W_A = 1.5\beta^{-1}$  for the forward transition and  $W_{\ddagger} - W_B = 1.8\beta^{-1}$  for the backward transition, respectively. For all four CLD models two Kramers's rates were obtained (cf. Sec. II), one via memory-free Kramers's theory (circles in Fig. 11) and the other by full inclusion of memory effects (stars in Fig. 11).

As can be seen from the figure, the transition rates of simulations  $CLD_{I-fit}$  and  $CLD_{II-fit}$  fall well into the range set by the reference trajectory (cf. horizontal lines), whereas the rates of  $CLD_I$  and  $CLD_{II}$  fall somewhat outside. The rates obtained with Kramers's theory did not differ significantly from the numerical results. Remarkably, all models yielded very similar rates with the memory-free and the full-memory version of Kramers's theory. This could suggest that memory effects do not influence transition rates significantly for the case at hand, and that integration of equations of motion could be simplified by replacing the generalized friction by a constant friction,  $\gamma_{eff} = \int_0^{\infty} \gamma(t) dt$ .

The transition rates obtained with constant friction, how-

ever, show that the opposite is true (cf. diamonds in Fig. 11). Integration with a constant friction significantly overestimates the rates for the models  $CLD_{I-fit}$  and  $CLD_{II-fit}$  and underestimates those obtained from  $CLD_I$  and  $CLD_{II}$ . Therefore, memory effects *do* play an important role.

It is somewhat surprising that the models  $CLD_I$  and  $CLD_{II}$  underestimated the transition rates despite the fact that their VACF is more accurate. Thus, despite the promise of method DIR to provide more accurate memory kernels than FIT, the opposite is the case. We attribute this to the numerical instabilities in solving the memory equation, which lead to several artifacts in the memory kernels  $\gamma_I$  and  $\gamma_{II}$  derived with DIR (cf. Sec. IV D). Accordingly, the stronger regularization applied in FIT avoids these artifacts that would otherwise compromise the rates. Further support is provided by the observation that damping of the artificially long tails of  $\gamma_I$  and  $\gamma_{II}$  led to improved transition rates (crosses in Fig. 11).

Overall, the results point to a sensitive influence of the 5 ps decay component on the transition rates: The artificially large component of DIR implies rates that are significantly lower than the reference. For TAIL, which leaves part of the artifact untouched, still a slight underestimation is observed. FIT, in contrast, which fully suppressed all slow components, overestimated the rate. Apparently, any further improvement would require a more accurate description of the slow component.

Effective friction constants provide a sensitive check even if the reference rate is unknown. Indeed, the effective friction  $\int_0^{\infty} \gamma(t) dt$  derived from the memory kernels  $\gamma_I$  and  $\gamma_{II}$  is too large ( $\approx 220 \text{ ps}^{-1}$ ), as can be seen by comparison with the estimate of an upper bound at  $120\text{--}150 \text{ ps}^{-1}$  obtained directly for the VACFs (Table II). In comparison, the memory kernels  $\gamma_{I-tail}$  and  $\gamma_{II-tail}$ , whose tail were damped down to zero (see Sec. IV D), yield effective friction constants of 114 and  $150 \text{ ps}^{-1}$ , respectively, which are at the upper boundary of the probable range of the true effective friction constants. Accordingly, they yielded improved transition rates (crosses in Fig. 11). Moreover, for the models  $CLD_{I-fit}$  and  $CLD_{II-fit}$  the transition rates were slightly too high, in agreement with the effective friction being underestimated due to the applied fit, which eradicates any slower

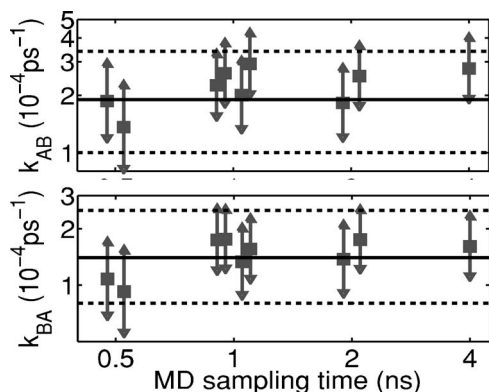


FIG. 12. Forward (top) and backward (bottom) transition rates predicted by fitted CLD models in dependence of the sampling time used to obtain the VACF from MD simulations. For the shortest sampling time of 0.5 ns only the rates obtained from memory kernels with the most extreme parameter values were computed. The horizontal lines depict the reference transition rate and its confidence interval. The boxes and their error bars depict predicted transition rates and their confidence intervals.

component of the VACF (cf. Table II). This suggests to use the friction integral as an additional and important regularization criterion for the memory kernel.

The CLD models discussed so far were based on VACFs obtained from simulations NT1 and NT2, with  $T > 9$  ns simulation time. To check if such a long simulation time is actually necessary to obtain sufficiently accurate memory kernels, we systematically assessed the amount of molecular dynamics sampling needed. To this end, eight 500 ps trajectories, NTS<sub>*i*</sub>,  $i = 1 \dots 8$ , were generated from different starting positions (cf. Sec. III) and used to compute memory kernels via the FIT method; for results see Table I. Memory kernels were computed from single trajectories, or from combinations of two, four, or all eight trajectories NTS<sub>*i*</sub>, constituting sampling times of 500 ps, 1 ns, 2 ns, and 4 ns, respectively. In Fig. 12 the transition rates predicted by the CLD model with these memory kernels are plotted against used sampling time. All obtained rates were within the range of and are centered at that of the reference MD simulation. The only exception are the rates obtained with memory kernels from the shortest sampling time (0.5 ns), which are systematically smaller than the reference MD rate (only the highest and lowest rates were shown). Already for sampling times  $t \geq 1$  ns the reference rates of the models CLD<sub>I-fit</sub> and CLD<sub>II-fit</sub> were reproduced. Thus, sampling as short as 1 ns is sufficient to correctly predict transition rates.

Taken together, the presented results show that the CLD model accurately predicts transition rates of complex systems. Remarkably, already the description of memory effects in terms of the VACF, which provides information mainly on short time correlations, proved sufficiently accurate to predict transition rates in NT at a time scale as long as 50 ns. Thus, not surprisingly sampling as short as 1 ns is sufficient, since the VACF is dominated by high frequency motion.

Furthermore, we found that the transition rate was mainly influenced by the effective friction  $\int_0^\infty \gamma(t) dt$ . Both methods to extract memory kernels, FIT and DIR, performed equally well, if the tailored memory kernels were used ( $\gamma_{I\text{-tail}}$ ,  $\gamma_{II\text{-tail}}$ ), whose effective friction was more accurate

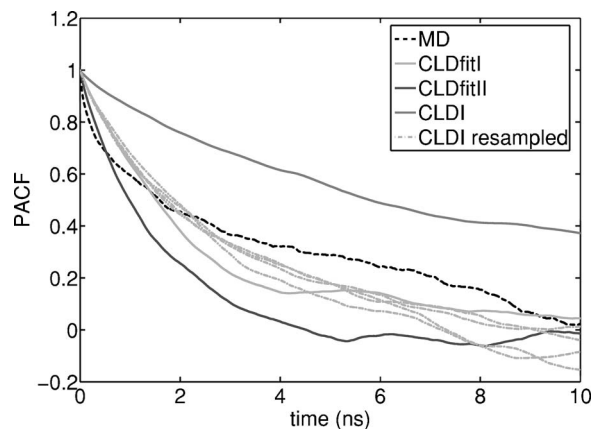


FIG. 13. Comparison of CLD generated (gray) PACFs with that computed from both MD trajectories NT1 and NT2. The curves of CLD<sub>I</sub> and CLD<sub>II-fit</sub> indicate the most extreme PACFs obtained from the discussed CLD models. The PACF of CLD<sub>I-fit</sub> agrees best with the MD results. From the scatter of the five PACFs for CLD<sub>I-fit</sub> (gray, dash dotted) the statistical error can be estimated.

than that of  $\gamma_I, \gamma_{II}$ . Thus, the reproduction of the medium scale oscillations of the VACF, which only the method DIR is capable of, was not important for the transition rate. However, an effective friction alone, i.e., a memory-free description, does not provide accurate rates (cf. diamonds in Fig. 11). Whereas memory effects were important on the short time scales of the integration steps, they were irrelevant on time scales probed by Kramers's theory. Indeed, already the memory obtained with FIT, decaying with ( $\tau \approx 1$  ps), influenced the dynamics over 100 integration steps, whereas the fastest time scale seen by Kramers's theory, i.e.,  $T \approx 2\pi/\omega = 9$  ps, is much slower and, therefore, is not affected by memory.

### G. Prediction of positional autocorrelation functions by CLD

The last observable of the CLD dynamics we compared to the reference MD is the positional autocorrelation function (PACF). Figure 13 shows the PACF obtained from MD simulations NT1 and NT2 covering a total simulation time of 153 ns and compared it to PACFs obtained from 300 ns CLD trajectories. We plotted the PACFs of models CLD<sub>I</sub> and CLD<sub>II-fit</sub> with slowest and fastest decays, respectively, as well as the PACF of CLD<sub>I-fit</sub>, which best agrees with the MD result.

The overall decay of all CLD-derived PACFs corresponds to that of the reference PACF from the MD simulation. Fits to single exponential decays yield decay times ranging from 135 to 9 ns for the CLD-derived PACFs, which are (roughly) on the same order of magnitude as that obtained from the MD-derived PACF (3.3 ns). Remarkably, the decay of the CLD-derived PACFs is systematically too slow for short times  $\tau < 0.5$  ns, whereas on long times some decays are faster and others slower than the reference. Moreover, the CLD-derived PACFs are well described by a single exponential decay, whereas the MD-derived PACF shows two significantly different time scales.

The large spread of the CLD-derived PACFs is striking.

In order to rule out that this is due to unconverged correlations we obtained several independent trajectories for each CLD model and computed their PACFs. For  $CLD_{I-fit}$ , these are shown in the figure, and their much smaller statistical variation confirms that the spread of the PACFs is indeed significant.

Furthermore, we compared the decay times of the CLD models with their respective transition rates. Correlation coefficients of  $r=0.69$  and  $r=0.72$  for the forward and backward rates, respectively, indicate a weak connection. However, the relatively low value also shows that not all dynamical properties that are relevant for the transition rates are captured by the PACF. Vice versa, other dynamical properties, which are described by the PACF, are not reflected in the transition rates.

The large differences between the PACFs are unexpected because they are uniquely defined by the corresponding VACFs, which vary much less for the different CLD models (cf. Fig. 10). Note, however, that the PACF is dominated by low frequency components, i.e., long time correlations, whereas the VACF is dominated by the high frequency components. The fact that the memory kernels computed from VACFs cannot capture these long time correlations explains the observed spread of the CLD-derived PACFs.

Nevertheless, the large spread of the PACFs indicates a tremendous influence of the memory kernels on the long time dynamics. In order to achieve better accuracy the PACF could be used in addition to the VACF to determine the memory kernel, e.g., solving the alternative memory equation, Eq. (20). However, here one needs to trade off the accuracy of the CLD model with the sampling time to gain the slowly converging PACF.

## V. CONCLUSIONS

Collective Langevin dynamics (CLD) provides a consistent framework to describe and simulate slow collective motions of proteins in an approach with a drastically reduced number of degrees of freedom and, hence, reduced dimensionality. In this framework the dynamics are separated into slow and fast degrees of freedom. The dynamics in the slow coordinates are evolved explicitly, whereas the fast degrees of freedom are treated in an implicit manner.

CLD is a bottom up approach based on first principles in the sense that all relevant information is extracted from the well validated description of protein dynamics by MD simulations. Furthermore, it is a systematic approach because the level of coarse graining can be tuned by the number of degrees of freedom which are explicitly considered. The extreme case of a one-dimensional description is presented here; the other extreme is explicit consideration of all degrees of freedom and in the CLD framework would trivially reproduce the MD model.

The slow nature of the conformational motion motivated and justified the application of the projection-operator formalism by Mori and Zwanzig to derive equations of motions for the dynamics of the collective coordinates. Both, linear (e.g., principal components) and curved coordinates were considered in full generality. The resulting exact equations of

motions take the form of a generalized Langevin equation with a potential of mean force. Here, we approximate this exact equation by replacing its *noise term* with a non-Markovian stochastic process that obeys the fluctuation dissipation theorem. The memory effects are found to be not negligible and thus are fully accounted for by a generalized frictional force, whose specific memory kernel is obtained for any dynamical system individually.

Memory kernels were computed from velocity autocorrelation functions obtained from short (few nanoseconds) MD trajectories via the corresponding Volterra-type equation. Because this inverse problem is notoriously difficult to solve and suffers from numerical instabilities, we tested different levels of regularization. The method FIT applied rather strong regularization, and hence was very robust against the inherent statistical noise in the VACF. In contrast, the second method, DIR, regularized only weakly, such that it allowed to capture more details of the VACF. The results indicated that for an accurate description of transition rates, the trade-off should be struck on the side of stronger regularization, i.e., increased robustness.

CLD is complementary and rests upon the many existing enhanced sampling methods to calculate free energy surfaces such as, REMD,<sup>43</sup> umbrella sampling,<sup>45,46</sup> or SMC.<sup>44</sup> All these methods, by construction, sacrifice dynamics to speed up sampling. We have proposed to reconstruct the conformational dynamics from the obtained free energy surfaces via CLD. Alternatively, ensembles obtained from experimental sources such as NMR might also be used to estimate a free energy surface.

As a test system, the hexapeptide neurotensin was considered. Explicit treatment in CLD was restricted to a one-dimensional (curved) conformational coordinate. Comparison of transition rates obtained from this extremely dimension reduced and, hence very efficient, description with those obtained from a 150 ns MD simulation showed an excellent agreement.

Remarkably, this good agreement for the neurotensin peptide was achieved by the most extreme conceivable dimension reduction, i.e., to only one dimension. A generalized curved coordinate was required to achieve such a drastic reduction; more than one but less than five linear degrees of freedom would likely allow to achieve similar accuracy.

We note that similar tests for much larger protein systems would of course be called for to further evaluate our approach. However, the requirement of converged reference transition rates from long MD simulations severely restricts the size of the test system. For instance, an available 450 ns simulation of crambin did not contain enough recurring transitions to reliably estimate reference transition rates, whereas enough transitions occurred in the presented 150 ns simulation of neurotensin. Nevertheless, our results indicate that CLD is also capable of accurately describing conformational dynamics of soluble proteins at microsecond time scales.

CLD yields trajectories with accurate thermodynamical and dynamical behaviors, in particular, accurate free energies and velocity autocorrelation functions. By focusing on relevant quantities, our CLD approach also provides new physical insights into the high-dimensional protein dynamics. The



relative fast decay of the memory kernel of neurotensin agrees with previous findings. For a similarly sized peptide an upper limit for a time scale on which no memory effect influenced transition rates was determined to be 1 ns.<sup>83</sup> This limit agrees with and is improved by our finding that memory effects did not play a significant role for transition rates at time scales above 10 ps. In focusing at accurate velocity autocorrelation functions, CLD might be particularly useful for the interpretation of neutron scattering experiments, which probe velocity autocorrelation functions.

The observed deviations of the CLD-derived *positional autocorrelation functions* indicate that for this observable memory effects on longer time scales are important. We further suggest to improve the accuracy of the required memory kernel by combining positional and velocity autocorrelation functions for its extraction, because the former probe long time scales and the latter short ones.

Whereas two or three explicit degrees of freedom can be treated within the presented framework in a straightforward manner, inclusion of more explicit coordinates will become impractical due to the high dimensionality of the free energy landscapes, which would render the nonparametric free energy estimation used here infeasible. As an alternative, weighted sums of multivariate Gaussians could be used to approximate the ensemble density. A CLD model based on a similar parametric approximation was already used here in the Kramers's approach, and its rates agreed well with those obtained from the nonparametric free energy surface.

We finally suggest that the extension to large conformational subspaces might allow on-the-fly computations of small regions of the free energy landscape, thereby, alleviating the sampling problem. In particular, the higher frequency PCA modes behave quasiharmonically and are much more efficiently sampled by MD than the low frequency modes. Thus, a two layered approach for CLD might be considered, which switches to explicit MD to probe entropic contributions to the free energy, whenever new, previously unknown, regions of the conformational subspace are encountered.

## ACKNOWLEDGMENTS

We thank Carsten Hartmann for helpful discussion on intricacies involved in a formulation of free energies for curved coordinates, Bert L. de Groot for providing the neurotensin trajectory NT1, and Rainer Böckmann for providing the HLA-B27 trajectory. We would also like to acknowledge the contribution by Andreas Briese, documented by his diploma thesis<sup>84</sup> supervised by one of us (H.G.) Our work has been supported by Volkswagenstiftung Grant Nos. I/80436 and I/78839.

## APPENDIX: DERIVATION OF A POTENTIAL OF MEAN FORCE

In this appendix we show that Eq. (15) evaluates to the force term  $\mathcal{P}\mathcal{L}|A_2(t)\rangle$ . To simplify notation, we use mass-weighted coordinates ( $\tilde{\mathbf{x}}=\mathbf{M}^{1/2}\mathbf{x}$  and  $\tilde{\mathbf{p}}=\mathbf{M}^{-1/2}\mathbf{p}$ ) and define

$$\delta_c := \delta_c,$$

$$\delta_v := \delta(\nabla_{\mathbf{x}}f \cdot \mathbf{p} - \dot{c}).$$

For the proof we will need the relations

$$\nabla_{\mathbf{x}}f \cdot \nabla_{\mathbf{x}}\delta_c = \delta'(f(\mathbf{x}) - c)\|\nabla_{\mathbf{x}}f\|^2, \quad (\text{A1})$$

$$\nabla_{\mathbf{x}}f \cdot \nabla_{\mathbf{x}}\delta_v = \nabla_{\mathbf{p}}\delta_v \cdot \nabla_{\mathbf{x}}(\nabla_{\mathbf{x}}f \cdot \mathbf{p}), \quad (\text{A2})$$

which are easily shown by applying the chain rule to the delta functions.

Consider the derivative of  $\rho(c, \dot{c})$ , which appears in Eq. (15). As seen from the definition, Eq. (13), its dependence on  $c$  is restricted to the delta-function  $\delta_c$ . Therefore,

$$\frac{\partial \rho(c, \dot{c})}{\partial c} = - \int \rho(\mathbf{x}, \mathbf{p}) \|\nabla_{\mathbf{x}}f\|^2 \delta'(f(\mathbf{x}) - c) \delta_v d\mathbf{x}d\mathbf{p},$$

which is transformed via relation Eq. (A1) to

$$\frac{\partial \rho(c, \dot{c})}{\partial c} = - \int \rho(\mathbf{x}, \mathbf{p}) \nabla_{\mathbf{x}}f \nabla_{\mathbf{x}}\delta_c \delta_v d\mathbf{x}d\mathbf{p}.$$

This allows us to eliminate the derivative of the delta-function  $\delta_c$  via partial integration,

$$\begin{aligned} \frac{\partial \rho(c, \dot{c})}{\partial c} = & \int \{ [(-\beta) \nabla_{\mathbf{x}}\mathcal{V} \cdot \nabla_{\mathbf{x}}f + \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}}f] \delta_v \\ & + \nabla_{\mathbf{x}}f \nabla_{\mathbf{x}}\delta_v \} \rho(\mathbf{x}, \mathbf{p}) \delta_c d\mathbf{x}d\mathbf{p}, \end{aligned}$$

where we have used that  $\rho(\mathbf{x}, \mathbf{p}) = Z^{-1} \exp(-\beta\mathcal{H})$ . To eliminate also the newly appeared derivative of  $\delta_v$ , Eq. (A2) is employed and the result integrated partially over momentum space. This yields

$$\begin{aligned} \frac{\partial \rho(c, \dot{c})}{\partial c} = & \int \{ -\beta \nabla_{\mathbf{x}}\mathcal{V} \cdot \nabla_{\mathbf{x}}f \\ & + \beta \nabla_{\mathbf{p}}\mathcal{H} \cdot \nabla_{\mathbf{x}}(\nabla_{\mathbf{x}}f \cdot \mathbf{p}) \} \rho(\mathbf{x}, \mathbf{p}) \delta_c d\mathbf{x}d\mathbf{p}, \end{aligned}$$

since the remaining terms

$$\nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}}f - \nabla_{\mathbf{p}} \cdot \nabla_{\mathbf{x}}(\nabla_{\mathbf{x}}f \cdot \mathbf{p})$$

cancel. Comparison with Eq. (14) shows that

$$\frac{1}{\rho(c, \dot{c})} \frac{\partial \rho(c, \dot{c})}{\partial c} = -\beta \mathcal{P}\mathcal{L}|A_2(t)\rangle.$$

<sup>1</sup>A. J. Wand, Nat. Struct. Biol. **8**, 926 (2001).

<sup>2</sup>J. Norberg and L. Nilsson, Q. Rev. Biophys. **36**, 257 (2003).

<sup>3</sup>R. H. Zhou, E. Harder, H. F. Xu, and B. J. Berne, J. Chem. Phys. **115**, 2348 (2001).

<sup>4</sup>M. E. Tuckerman, B. J. Berne, and G. J. Martyna, J. Chem. Phys. **94**, 6811 (1991).

<sup>5</sup>M. Tuckerman, B. J. Berne, and G. J. Martyna, J. Chem. Phys. **97**, 1990 (1992).

<sup>6</sup>P. Minary, M. E. Tuckerman, and G. J. Martyna, Phys. Rev. Lett. **93**, 150201 (2004).

<sup>7</sup>J. L. Scully and J. Hermans, Mol. Simul. **11**, 67 (1993).

<sup>8</sup>F. Zhang, J. Chem. Phys. **106**, 6102 (1997).

<sup>9</sup>J. A. Board, J. W. Csey, J. F. Leathrum, A. Windemuth, and K. Schulten, Chem. Phys. Lett. **198**, 89 (1992).

<sup>10</sup>A. M. Mathiowetz, A. Jain, N. Karasawa, and W. A. Goddard, Proteins **20**, 227 (1994).

<sup>11</sup>M. Eichinger, H. Grubmüller, H. Heller, and P. Tavan, J. Comput. Chem. **18**, 1729 (1997).

<sup>12</sup>L. Greengard and V. Rokhlin, Chem. Scr. **29A**, 139 (1989).

<sup>13</sup>A. Y. Toukumaji and J. A. Board, Comput. Phys. Commun. **95**, 73 (1996).

- <sup>14</sup> J. P. Ryckaert, G. Ciccotti, and H. J. C. Berendsen, *J. Comput. Phys.* **23**, 327 (1977).
- <sup>15</sup> S. Miyamoto and P. A. Kollman, *J. Comput. Chem.* **13**, 952 (1992).
- <sup>16</sup> B. Hess, H. Bekker, H. J. C. Berendsen, and J. G. E. M. Fraaije, *J. Comput. Chem.* **18**, 1463 (1997).
- <sup>17</sup> K. Tai, *Biophys. Chem.* **107**, 213 (2004).
- <sup>18</sup> N. Go and H. A. Scheraga, *J. Chem. Phys.* **51**, 4751 (1969).
- <sup>19</sup> B. Roux and T. Simonson, *Biophys. Chem.* **78**, 1 (1999).
- <sup>20</sup> T. Head-Gordon and S. Brown, *Curr. Opin. Struct. Biol.* **13**, 160 (2003).
- <sup>21</sup> S. J. Marrink and D. P. Tieleman, *Biophys. J.* **83**, 2386 (2002).
- <sup>22</sup> G. Ayton and G. A. Voth, *Biophys. J.* **83**, 3357 (2002).
- <sup>23</sup> J. P. Ulmschneider and W. L. Jorgensen, *J. Chem. Phys.* **118**, 4261 (2003).
- <sup>24</sup> F. Sartori, B. Melchers, H. Bottcher, and E. W. Knapp, *J. Chem. Phys.* **108**, 8264 (1998).
- <sup>25</sup> A. Kloczkowski, J. E. Mark, and B. Erman, *Macromolecules* **22**, 1423 (1989).
- <sup>26</sup> A. Amadei, A. B. M. Linssen, and H. J. C. Berendsen, *Proteins* **17**, 412 (1993).
- <sup>27</sup> A. E. Garcia, *Phys. Rev. Lett.* **68**, 2696 (1992).
- <sup>28</sup> A. Amadei, A. B. M. Linssen, B. L. de Groot, D. M. F. van Aalten, and H. J. C. Berendsen, *J. Biomol. Struct. Dyn.* **13**, 615 (1996).
- <sup>29</sup> R. Zwanzig, *Phys. Rev.* **124**, 983 (1961).
- <sup>30</sup> H. Mori, *Prog. Theor. Phys.* **33**, 423 (1965).
- <sup>31</sup> I. Benjamin, L. L. Lee, Y. S. Li, A. Liu, and K. R. Wilson, *Chem. Phys.* **152**, 1 (1991).
- <sup>32</sup> R. Ferrando, R. Spadacini, and G. E. Tommei, *Chem. Phys. Lett.* **347**, 487 (2001).
- <sup>33</sup> C. C. Martens, *J. Chem. Phys.* **116**, 2516 (2002).
- <sup>34</sup> W. K. Park and S. C. Park, *J. Mol. Struct.: THEOCHEM* **630**, 215 (2003).
- <sup>35</sup> P. Romiszowski and R. Yaris, *J. Chem. Phys.* **94**, 6751 (1991).
- <sup>36</sup> Y. X. Zhang and S. C. Park, *Bull. Korean Chem. Soc.* **21**, 1095 (2000).
- <sup>37</sup> M. A. Balsera, W. Wriggers, Y. Oono, and K. Schulten, *J. Phys. Chem.* **100**, 2567 (1996).
- <sup>38</sup> A. Amadei, M. A. Ceruso, and A. DiNola, *Proteins* **36**, 419 (1999).
- <sup>39</sup> B. J. Berne and G. D. Harp, *Adv. Chem. Phys.* **17**, 63 (1970).
- <sup>40</sup> G. R. Kneller and K. Hinsen, *J. Chem. Phys.* **115**, 11097 (2001).
- <sup>41</sup> J. P. Boon and S. A. Rice, *J. Chem. Phys.* **47**, 2480 (1967).
- <sup>42</sup> P. A. Egelstaff, *Phys. Chem. Liq.* **16**, 293 (1987).
- <sup>43</sup> Y. Sugita and Y. Okamoto, *Chem. Phys. Lett.* **314**, 141 (1999).
- <sup>44</sup> A. Kidera, *Int. J. Quantum Chem.* **75**, 207 (1999).
- <sup>45</sup> G. M. Torrie and J. P. Valle, *J. Chem. Phys.* **66**, 1402 (1977).
- <sup>46</sup> M. Souaille and B. Roux, *Comput. Phys. Commun.* **135**, 40 (2001).
- <sup>47</sup> H. Grubmüller and P. Tavan, *J. Comput. Chem.* **19**, 1534 (1998).
- <sup>48</sup> P. Hänggi, P. Talkner, and M. Borkovec, *Rev. Mod. Phys.* **62**, 251 (1990).
- <sup>49</sup> H. Heise, S. Luca, B. L. de Groot, H. Grubmüller, and M. Baldus, *Biophys. J.* **89**, 2113 (2005).
- <sup>50</sup> B. J. Berne, *Dynamic Light Scattering* (Wiley, New York, 1976).
- <sup>51</sup> R. Zwanzig, in *Systems Far from Equilibrium*, edited by L. Garrido (Springer, New York, 1980), pp. 198–225.
- <sup>52</sup> R. Kubo, *Rep. Prog. Phys.* **29**, 255 (1966).
- <sup>53</sup> A. J. Chorin, O. H. Hald, and R. Kupferman, *Proc. Natl. Acad. Sci. U.S.A.* **97**, 2968 (2000).
- <sup>54</sup> G. D. Harp and B. J. Berne, *Phys. Rev. A* **2**, 975 (1970).
- <sup>55</sup> H. Grubmüller, *Phys. Rev. E* **52**, 2893 (1995).
- <sup>56</sup> M. Berkowitz, J. D. Morgan, and J. A. McCammon, *J. Chem. Phys.* **78**, 3256 (1983).
- <sup>57</sup> T. Srokowski, *Phys. Rev. E* **64**, 031102 (2001).
- <sup>58</sup> M. Iannuzzi, A. Laio, and M. Parrinello, *Phys. Rev. Lett.* **90**, 238302 (2003).
- <sup>59</sup> J. P. Boon and S. Yip, *Molecular Hydrodynamics* (McGraw-Hill, New York, 1980).
- <sup>60</sup> F. Shimojo, K. Hoshino, and M. Watabe, *J. Phys. Soc. Jpn.* **63**, 141 (1994).
- <sup>61</sup> M. Berkowitz, J. D. Morgan, D. J. Kouri, and J. A. McCammon, *J. Chem. Phys.* **75**, 2462 (1981).
- <sup>62</sup> D. E. Smith and C. B. Harris, *J. Chem. Phys.* **92**, 1304 (1990).
- <sup>63</sup> T. Shimizu, *Physica A* **164**, 123 (1990).
- <sup>64</sup> E. Lindahl, B. Hess, and D. Van der Spoel, *J. Mol. Model.* **7**, 306 (2001); <http://www.gromacs.org>
- <sup>65</sup> T. Darden, D. York, and L. Pedersen, *J. Chem. Phys.* **98**, 10089 (1993).
- <sup>66</sup> U. Essmann, L. Perara, M. L. Berkowitz, T. Darda, H. Lee, and L. G. Pedersen, *J. Chem. Phys.* **103**, 8577 (1995).
- <sup>67</sup> H. J. C. Berendsen, J. P. M. Postma, A. DiNola, and J. R. Haak, *J. Chem. Phys.* **81**, 3684 (1984).
- <sup>68</sup> W. L. Jorgensen, D. S. Maxwell, and J. Tirado-Rives, *J. Am. Chem. Soc.* **118**, 11225 (1996).
- <sup>69</sup> C. T. H. Baker, *J. Comput. Appl. Math.* **125**, 217 (2000).
- <sup>70</sup> P. K. Lamm, *Surveys on Solution Methods for Inverse Problems* (Springer, New York, 2000), pp. 53–82.
- <sup>71</sup> P. K. Lamm and L. Elden, *SIAM (Soc. Ind. Appl. Math.) J. Numer. Anal.* **34**, 1432 (1997).
- <sup>72</sup> O. Lange, Ph.D. thesis, Georg-Augustus-Universität, 2005.
- <sup>73</sup> P. C. Hansen, *SIAM Rev.* **34**, 561 (1992).
- <sup>74</sup> P. C. Hansen and D. P. O'Leary, *SIAM J. Sci. Comput. (USA)* **14**, 1487 (1993).
- <sup>75</sup> H. W. Engl, M. Hanke, and A. Neubauer, *Regularization of Inverse Problems* (Kluwer, Dordrecht, 2000).
- <sup>76</sup> M. P. Allen and D. J. Tildesley, *Computer Simulation of Liquids* (Clarendon, Oxford, 1987).
- <sup>77</sup> M. E. Tuckerman and B. J. Berne, *J. Chem. Phys.* **95**, 4389 (1991).
- <sup>78</sup> J. W. Cooley and J. W. Tukey, *Math. Comput.* **19**, 297 (1965).
- <sup>79</sup> H. J. Nussbaumer, *Fast Fourier Transform and Convolution Algorithms* (Springer, New York, 1982).
- <sup>80</sup> C. T. H. Baker and M. S. Derakhshan, *J. Comput. Appl. Math.* **20**, 5 (1987).
- <sup>81</sup> C. Schütte, A. Fischer, W. Huisinga, and P. Deuffhard, *J. Comput. Phys.* **151**, 146 (1999).
- <sup>82</sup> H. Frauenfelder and P. G. Wolynes, *Science* **229**, 337 (1985).
- <sup>83</sup> B. L. de Groot, X. Daura, A. E. Mark, and H. Grubmüller, *J. Mol. Biol.* **309**, 299 (2001).
- <sup>84</sup> A. Briese, Master's thesis, Ludwig-Maximilians-Universität, 1996.
- <sup>85</sup> O. Lange, Ph.D. thesis, "Collective Langevin Dynamics of Conformational Motions in Proteins," Cuvillier Verlag Gröttingen, 2006.