# Molecular dynamics simulation of polymer helix formation using rigid-link methods

D. C. Rapaport*

*Physics Department, Bar-Ilan University, Ramat-Gan 52900, Israel*
(Received 14 February 2002; published 15 July 2002)

Molecular dynamics simulations are used to study structure formation in simple model polymer chains that are subject to excluded volume and torsional interactions. The changing conformations exhibited by chains of different lengths under gradual cooling are followed until each reaches a state from which no further change is possible. The interactions are chosen so that the true ground state is a helix, and a high proportion of simulation runs succeed in reaching this state; the fraction that manages to form defect-free helices is a function of both chain length and cooling rate. In order to demonstrate behavior analogous to the formation of protein tertiary structure, additional attractive interactions are introduced into the model, leading to the appearance of aligned, antiparallel helix pairs. The simulations employ a computational approach that deals directly with the internal coordinates in a recursive manner; this representation is able to maintain constant bond lengths and angles without the necessity of treating them as an algebraic constraint problem supplementary to the equations of motion.

## I. INTRODUCTION

Polymers, because of their importance and complexity, have provided a longstanding challenge for computer simulation. Over the years, the field has become fragmented, both in terms of the problems addressed and the methodology employed. Broadly speaking, the kinds of system studied can be classified into distinct groups; there are biological heteropolymers, a category dominated by the proteins; homopolymers and block copolymers that include a great variety of molecular types, from alkanes to plastics; and idealized polymer models used for elucidating general principles such as the theta point, reptation, and multiphase behavior. The computational techniques span an equally broad range; they include molecular dynamics (MD) simulation employing models that represent the molecules at various levels of detail, ranging from fully atomic to highly reduced descriptions; Monte Carlo sampling of both continuum- and lattice-based systems, again with different levels of representation; and exact enumeration of small systems aimed at eliminating the sampling errors inherent in the other methods. While all three kinds of methodology provide important information about equilibrium behavior and, in a sense, amount to doing statistical mechanics numerically, the MD approach provides access to the dynamical and nonequilibrium aspects of the behavior; although it might be argued that Monte Carlo shares some of this capability, the associated dynamics is a consequence of the chosen stochastic sampling algorithm. Lattice-based approaches, though offering a vastly reduced configuration space, have the additional problem of the discreteness of the lattice on which the polymer is embedded, and the consequent absence of gradual transitions between different configurations.

The inherent difficulty in polymer simulation is that the problem naturally embraces a broad range of time scales, ranging from very fast processes associated with bond vibra-

tion, followed by the somewhat slower, highly localized conformational changes such as crankshaft motions, then the even slower aspects of reorganization such as the still relatively localized process of helix formation, and, finally, the typically extremely slow changes that lead to the emergence of tertiary structure characteristic of protein folding and to polymer diffusion in a concentrated solution. The time scales associated with this hierarchy of processes span a range considerably in excess of ten orders of magnitude, and so such systems are clearly not generally amenable to direct modeling, unless subjected to major simplification. Considerable effort has been invested in the design of models and simulation methods with the aim of alleviating this problem to at least some degree.

One especially important application of polymer simulation is in the field of protein folding, e.g., Refs. [1–5]; achieving an understanding of the mechanisms underlying this important process presents a major challenge to computational biochemistry. Protein modeling runs the gamut from, at one extreme, highly detailed molecular representations involving potentials derived from a mixture of theory and experiment, together with a solvent of individual water molecules, all solved by MD and an enormous amount of computational effort [6,7], through highly simplified models also solved by MD [8], to yet even simpler models embedded in lattices with only a limited number of degrees of freedom (DOFs) studied using a suitable Monte Carlo procedure and a greatly reduced investment in computing [2]; even complete enumeration of all conformations is sometimes feasible [9]. While the manner in which the amino acid sequence of any given protein is able to determine its presumably unique spatial structure continues to be the subject of intense study, of no less importance is the question of the folding pathway—the preferred route (or routes) through multidimensional conformation space eventually terminating at the native state. While all the widely differing methodologies enumerated above can be used for studying folded states, the collective dynamical processes that underlie folding really demand an approach based on MD. But, after

---

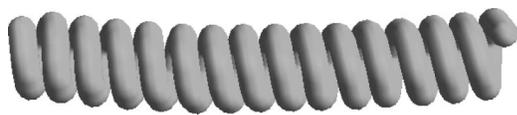*Electronic address: rapaport@mail.biu.ac.il

FIG. 1. A well-formed helix in a chain of length 90; a goal of the simulations is to observe chains spontaneously collapsing into this state (the polymer is drawn as a tube whose radius is that of the monomers).

reaching this conclusion, there is a practical question of whether, even after substantial simplification, serious progress in understanding the mechanisms of folding can be achieved by computer simulation, owing to the diversity of intrinsic time scales; while substantial advances have been made, a great deal remains to be done before this question is answered.

The goal of the present paper is twofold. The first goal is a demonstration of a different perspective on the MD approach to studying protein folding. The most ambitious level of modeling is based on carefully constructed potential functions, often with a multitude of parameters; since the native conformation generally corresponds to the state of minimum free energy, establishing the details of these interatomic interactions, including solvent effects, provides the foundation for such work. Determining whether the known native state of a given protein is the one favored by energetic considerations is in itself a complex optimization task, but following the full dynamics over a sufficiently long period of time for the major structural changes that typify protein folding to occur verges on the impossible. The approach adopted here is just the opposite, and the question posed is the following. Given a known structural motif, such as the helix, and a simplified model of a polymer chain with a readily determined, unique ground state corresponding to this configuration, as in Fig. 1, will the chain collapse into this state within a reasonable amount of computation time when allowed to move freely in space, as shown in Fig. 2, while subjected to gradual cooling?

The most elementary of these organized structures is the helix, which, while being a prominent feature in many globular proteins, is only classified as a secondary structural element (the primary structure being the amino acid sequence itself), and because of its homogeneous nature (except for the
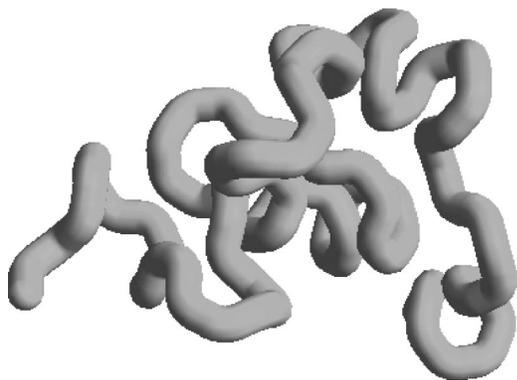
ends) it might be argued that being able to fold a helix is not really a significant step in learning how to fold an entire protein. Therefore, another folding problem considered here is the one with a ground state formed from an antiparallel pair of helices. This, too, is a recognizable element in some proteins, and is unquestionably classified as a tertiary structure.

The obvious extension of this approach, a subject for future exploration, is to design simple models for other structural motifs, with the hope of learning more about folding by examining the collapse pathways of these idealized models; some structures might fold more readily than others, in which case the steric and topological issues involved could be investigated; for some structures there might be recognizable intermediate states along the folding trajectory; some cases might reveal useful properties that, when regarded as conformational (or reaction) ''coordinates,'' might serve in the design of other kinds of simplified models [9]; and finally, once the simple version has been found to have the correct behavior, the models could be enhanced by gradually incorporating features from more realistic representations, including specific interactions and structural details. This represents the motivation for this kind of modeling approach.

The second goal is methodological. Even when considering the simplest of model polymers, in which, typically, all the molecular detail is absorbed into effective atoms located along the backbone chain (more so if this simplification is not made) there is a need to specify the internal DOFs of the system. One possibility is to assume that adjacent atoms are connected by stiff springs represented by a suitable potential function; in this case each atom has its full complement of three translational DOFs and, if these atoms are regarded as rigid particles rather than point masses, three rotational DOFs as well. If the bond potentials are made sufficiently stiff to correspond to a typical real system, the ensuing high-frequency vibrations impose a very small integration time step, which runs contrary to the goal of efficiently simulating over long periods of time.

It is, however, possible to introduce geometrical restrictions, such as strictly constant bond lengths, while retaining a soluble dynamical problem. This is done by introducing holonomic constraints and Lagrange multipliers into the equations of motion [10], and then solving a set of algebraic equations while integrating the differential equations of motion. Two approaches have been developed for doing this; one involves initially solving the unconstrained equations of motion over a single time step and then iteratively correcting the relative coordinates [11,12], and, optionally, also the relative velocities [13], using a relaxation procedure to ensure the constraints remain satisfied; the other tackles the problem by constructing a matrix representing the contributions of the constraints which, in effect, must be inverted at each time step [14,15], and which is subject to gradual drift requiring regular correction. Similar geometric constraints can be introduced to maintain constant bond angles as well, since it is often a reasonable approximation to assume that the angles between consecutive bonds along the backbone (or elsewhere) are unvarying. Such geometrical constraints have proved extremely useful, given the nature of the excitations



FIG. 2. A randomly coiled chain of length 90; this configuration represents a typical state of the chain prior to the onset of folding.

present in the system: fluctuations in bond lengths, and sometimes also angles, tend to be of relatively small amplitude and high frequency, so that freezing them out of the dynamics permits a substantial increase in the allowed integration time step. The amount of additional processing required for the constraints depends on their number $n_c$; the dependence is typically $O(n_c)$ for the iterative approach, but for the matrix approach it is $O(n_c^3)$, making the latter unsuitable for large problems.

If bonds lengths and angles are fixed, the only remaining internal DOFs are the dihedral angles, each defined in terms of a rotation about an axis lying along a bond, and affecting the relative orientation of the pair of bonds on either side. For reasons shrouded in history, dealing with this problem has been perceived as difficult, as indeed it is, if the problem is not addressed in a suitable manner. A significant advance in the methodology for dealing with dynamical problems involving internal coordinates occurred some years ago in the robotics field [16,17], but with only the occasional exception, e.g., Ref. [18], it appears to have gone unappreciated by the polymer simulation community at large. Because of the importance of this technique, the goal of which is to deal directly and economically with the internal DOFs, and since there is no reason why it should not be capable of replacing the various constraint-based approaches for most applications, a detailed treatment of the underlying theory is included in the paper.

This approach to the dynamics of linked bodies also requires solving the dynamics of individual rigid bodies. An alternative, recently described means [19] of numerically dealing with the rigid-body equations of motion is discussed briefly; the method is based on rotation matrices, rather than on quaternions (or even Euler angles) that are generally used. The present formulation differs slightly from the original in regard to the reference frame in which the computations are carried out. The use of rotation matrices offers improved numerical stability, and since the method belongs to the leapfrog family of integrators, it means that simple leapfrog integration techniques can be used for the entire set of dynamical equations appearing in the problem.

## II. LINKED-BODY DYNAMICS

### A. Chain coordinates

Consider a linear polymer chain whose monomers are joined by rigid bonds. In the discussion that follows, the terms "monomer," "atom," "site," and "joint" will be used interchangeably, as appropriate to the context, likewise "link" and "bond." Bond lengths and angles are constant. If each torsional DOF is regarded as a mechanical joint associated with the site at one end of the link, with just a single rotational DOF, then the system is analogous to a basic problem in the field of robotic manipulators [16,17].

The chain configuration is defined by the site positions $\{r_k\}$, and if the bond vectors between adjacent sites are $\{b_k\}$ then $r_{k+1}=r_k+b_k$. The internal configuration of the chain can be specified by a set of bond rotation matrices $\{R_k\}$. The transformation between the local coordinate frames attached to bonds $k-1$ and $k(k \geqslant 1)$ involves a rotation through the

bond angle $\alpha_k$ about the axis $\hat{x}_{k-1}$, where $\cos \alpha_k = \hat{b}_{k-1} \cdot \hat{b}_k$, followed by a rotation through the dihedral angle $\theta_k$ about the joint axis $\hat{z}_{k-1}$. The matrix (actually its transpose) corresponding to this rotation is

$$R_{k-1,k}^T = \begin{pmatrix} \cos \theta_k & -\sin \theta_k \cos \alpha_k & \sin \theta_k \sin \alpha_k \\ \sin \theta_k & \cos \theta_k \cos \alpha_k & -\cos \theta_k \sin \alpha_k \\ 0 & \sin \alpha_k & \cos \alpha_k \end{pmatrix},$$

$$(1)$$

so that

$$R_k^T = R_0^T R_{0,1}^T \cdots R_{k-1,k}^T, \tag{2}$$

where $R_0^T$ represents the orientation of the initial site and bond, and

$$r_{k+1} = r_k + |b_k| R_k^T \hat{z}. \tag{3}$$

In the present case, $\{|b_k|\}$ and $\{\alpha_k\}$ are all constant, so that the only internal DOFs are those associated with $\{\theta_k\}$. Define $\hat{h}_k$ to be the rotation axis of the joint between bonds $k-1$ and $k$ that is fixed in the frame of bond $k-1$; in the present case $\hat{h}_k \equiv \hat{z}_{k-1}$. Insofar as indexing is concerned, there are $n_r$ internal rotational joints (with labels $1, \ldots, n_r$), $n_b = n_r + 1$ bonds $(0, \ldots, n_r)$, and $n_r + 2$ sites $(0, \ldots, n_r + 1)$. In order to completely specify the chain configuration, an additional joint is attached to the $k=0$ site, with three translational and three rotational DOFs (conceptually equivalent to a telescopic ball-and-socket joint); this joint is included in the formalism but will, eventually, be treated separately.

### B. Kinematic and dynamic relations

If $v_k$ and $\omega_k$ are the linear and angular velocities of site $k$, then the velocities and accelerations of adjacent sites are related by

$$\omega_k = \omega_{k-1} + \hat{h}_k \dot{\theta}_k, \tag{4}$$

$$v_k = v_{k-1} + \omega_{k-1} \times b_{k-1}, \tag{5}$$

$$\dot{\omega}_k = \dot{\omega}_{k-1} + \hat{h}_k \ddot{\theta}_k + \omega_{k-1} \times \hat{h}_k \dot{\theta}_k, \tag{6}$$

$$\dot{v}_k = \dot{v}_{k-1} + \dot{\omega}_{k-1} \times b_{k-1} + \omega_{k-1} \times (\omega_{k-1} \times b_{k-1}), \tag{7}$$

where $1 \leqslant k \leqslant n_r$. While the mass elements of the chain are normally identified with the sites, here it is helpful to associate them with the bonds; if $r_k + c_k$ is the location of the center of mass of the atoms attached to bond $k$, then the center-of-mass acceleration of the bond is

$$\dot{v}_k^c = \dot{v}_k + \dot{\omega}_k \times c_k + \omega_k \times (\omega_k \times c_k). \tag{8}$$

If $f_k$ and $n_k$ are the force and torque acting on bond $k$ across joint $k$, then the equations of motion are

$$\mathcal{I}_k \dot{\boldsymbol{\omega}}_k + \boldsymbol{\omega}_k \times (\mathcal{I}_k \boldsymbol{\omega}_k) = \boldsymbol{n}_k - \boldsymbol{n}_{k+1} - \boldsymbol{c}_k \times \boldsymbol{f}_k$$
$$- (\boldsymbol{b}_k - \boldsymbol{c}_k) \times \boldsymbol{f}_{k+1} + \boldsymbol{n}_k^e , \qquad (9)$$

$$m_k \dot{\boldsymbol{v}}_k^c = \boldsymbol{f}_k - \boldsymbol{f}_{k+1} + \boldsymbol{f}_k^e , \qquad (10)$$

where $\boldsymbol{f}_k^e$ and $\boldsymbol{n}_k^e$ are the externally applied force and torque; $m_k$ and $\mathcal{I}_k$ are the mass and moment of inertia of (the atoms associated with) the bond, the latter expressed in a space-fixed frame and relative to the center of mass of the bond. It is often convenient when dealing with rigid bodies to work in a center-of-mass frame [10]; this is not the case here, and all vector components are expressed in the space-fixed coordinate frame. Rearrange the terms of Eqs. (9) and (10) to obtain relations between torques and forces on adjacent bonds,

$$\boldsymbol{n}_k = \boldsymbol{n}_{k+1} + \boldsymbol{b}_k \times \boldsymbol{f}_{k+1} + m_k \boldsymbol{c}_k \times \dot{\boldsymbol{v}}_k^c + \mathcal{I}_k \dot{\boldsymbol{\omega}}_k + \boldsymbol{\omega}_k \times (\mathcal{I}_k \boldsymbol{\omega}_k) - \boldsymbol{n}_k^e$$
$$- \boldsymbol{c}_k \times \boldsymbol{f}_k^e , \qquad (11)$$

$$\boldsymbol{f}_k = \boldsymbol{f}_{k+1} + m_k \dot{\boldsymbol{v}}_k^c - \boldsymbol{f}_k^e , \qquad (12)$$

and define the torque

$$t_k = \hat{\boldsymbol{h}}_k \cdot \boldsymbol{n}_k \qquad (13)$$

that acts along the axis $\hat{\boldsymbol{h}}_k$ at joint $k$ and corresponds to the torsional interaction due to a twist around bond $k-1$.

### C. Spatial operator formulation

Equations (4)–(7) can be expressed more concisely in terms of six-component "spatial" vectors that combine the translational and rotational quantities. It is also convenient to represent certain vectors by means of antisymmetric matrices of form

$$\tilde{\boldsymbol{u}} = \begin{pmatrix} 0 & -u_z & u_y \\ u_z & 0 & -u_x \\ -u_y & u_x & 0 \end{pmatrix}, \qquad (14)$$

so that $\tilde{\boldsymbol{u}} \boldsymbol{v} \equiv \boldsymbol{u} \times \boldsymbol{v}$. The resulting equations are

$$\begin{pmatrix} \boldsymbol{\omega}_k \\ \boldsymbol{v}_k \end{pmatrix} = \begin{pmatrix} I & 0 \\ -\tilde{\boldsymbol{b}}_{k-1} & I \end{pmatrix} \begin{pmatrix} \boldsymbol{\omega}_{k-1} \\ \boldsymbol{v}_{k-1} \end{pmatrix} + \begin{pmatrix} \hat{\boldsymbol{h}}_k \\ 0 \end{pmatrix} \dot{\theta}_k , \qquad (15)$$

$$\begin{pmatrix} \dot{\boldsymbol{\omega}}_k \\ \dot{\boldsymbol{v}}_k \end{pmatrix} = \begin{pmatrix} I & 0 \\ -\tilde{\boldsymbol{b}}_{k-1} & I \end{pmatrix} \begin{pmatrix} \dot{\boldsymbol{\omega}}_{k-1} \\ \dot{\boldsymbol{v}}_{k-1} \end{pmatrix} + \begin{pmatrix} \hat{\boldsymbol{h}}_k \\ 0 \end{pmatrix} \ddot{\theta}_k$$
$$+ \begin{pmatrix} \boldsymbol{\omega}_{k-1} \times \hat{\boldsymbol{h}}_k \dot{\theta}_k \\ \boldsymbol{\omega}_{k-1} \times (\boldsymbol{\omega}_{k-1} \times \boldsymbol{b}_{k-1}) \end{pmatrix}, \qquad (16)$$

or, equivalently,

$$V_k = \phi_{k-1,k}^T V_{k-1} + H_k^T \dot{W}_k , \qquad (17)$$

$$A_k = \phi_{k-1,k}^T A_{k-1} + H_k^T \ddot{W}_k + X_k , \qquad (18)$$

where $V_k$ and $A_k$ are examples of spatial vectors, and

$$\phi_{k-1,k}^T = \begin{pmatrix} I & 0 \\ -\tilde{\boldsymbol{b}}_{k-1} & I \end{pmatrix}. \qquad (19)$$

The $6 \times 6$ matrices $\phi_{k-1,k}^T$ and $\phi_{k,k+1}$ (later) appear throughout the derivation, and their role is to propagate kinematic and dynamic information between joints. Several other new variables have been used,

$$H_k^T = \begin{pmatrix} \hat{\boldsymbol{h}}_k \\ 0 \end{pmatrix} \qquad (20)$$

is a six-component joint axis vector (in the more general case of a joint with $d$ DOFs, which the formalism is capable of handling, $H_k^T$ would become a $6 \times d$ matrix),

$$X_k = \begin{pmatrix} \tilde{\boldsymbol{\omega}}_{k-1} & 0 \\ 0 & \tilde{\boldsymbol{\omega}}_{k-1} \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{h}}_k \dot{\theta}_k \\ \boldsymbol{v}_k - \boldsymbol{v}_{k-1} \end{pmatrix} \qquad (21)$$

is a six-component spatial vector containing the remaining acceleration terms of the current site, and $\dot{W}_k \equiv \dot{\theta}_k$. When used in vectors and matrices, $I$ and $0$ denote unit and zero block submatrices of the implied size. The six-component vectors, and most of the associated matrices, are shown in capital italic letters (to retain some similarity with Ref. [20], $\phi$, $\psi$, and $\mathcal{M}$ are also used); no other special notation is needed since the variable types will be obvious from the context.

In a similar way, Eqs. (11)–(13) can be rewritten as

$$\begin{pmatrix} \boldsymbol{n}_k \\ \boldsymbol{f}_k \end{pmatrix} = \begin{pmatrix} I & \tilde{\boldsymbol{b}}_k \\ 0 & I \end{pmatrix} \begin{pmatrix} \boldsymbol{n}_{k+1} \\ \boldsymbol{f}_{k+1} \end{pmatrix}$$
$$+ \begin{pmatrix} m_k \boldsymbol{c}_k \times \dot{\boldsymbol{v}}_k^c + \mathcal{I}_k \dot{\boldsymbol{\omega}}_k + \boldsymbol{\omega}_k \times (\mathcal{I}_k \boldsymbol{\omega}_k) \\ m_k \dot{\boldsymbol{v}}_k^c \end{pmatrix}$$
$$- \begin{pmatrix} \boldsymbol{n}_k^e + \boldsymbol{c}_k \times \boldsymbol{f}_k^e \\ \boldsymbol{f}_k^e \end{pmatrix}, \qquad (22)$$

$$\begin{pmatrix} t_k \\ 0 \end{pmatrix} = \begin{pmatrix} \hat{\boldsymbol{h}}_k \\ 0 \end{pmatrix}^T \begin{pmatrix} \boldsymbol{n}_k \\ \boldsymbol{f}_k \end{pmatrix}, \qquad (23)$$

or, equivalently,

$$F_k = \phi_{k,k+1} F_{k+1} + M_k A_k + Y_k , \qquad (24)$$

$$T_k = H_k F_k . \qquad (25)$$

Here Eq. (8) has been used, and

$$M_k = \begin{pmatrix} \mathcal{I}_k - m_k \tilde{\boldsymbol{c}}_k \tilde{\boldsymbol{c}}_k & m_k \tilde{\boldsymbol{c}}_k \\ -m_k \tilde{\boldsymbol{c}}_k & m_k I \end{pmatrix}, \qquad (26)$$

is the symmetric, $6 \times 6$ mass matrix; the six-component vector

$$Y_k = \begin{pmatrix} \tilde{\boldsymbol{\omega}}_k (\mathcal{I}_k - m_k \tilde{\boldsymbol{c}}_k \tilde{\boldsymbol{c}}_k) \boldsymbol{\omega}_k \\ m_k \tilde{\boldsymbol{\omega}}_k \tilde{\boldsymbol{\omega}}_k \boldsymbol{c}_k \end{pmatrix} - \begin{pmatrix} \boldsymbol{n}_k^e + \boldsymbol{c}_k \times \boldsymbol{f}_k^e \\ \boldsymbol{f}_k^e \end{pmatrix} \quad (27)$$

contains the remaining force contributions. The identity $\boldsymbol{c}_k \times [\boldsymbol{\omega}_k \times (\boldsymbol{\omega}_k \times \boldsymbol{c}_k)] = -\boldsymbol{\omega}_k \times [\boldsymbol{c}_k \times (\boldsymbol{c}_k \times \boldsymbol{\omega}_k)]$ was used in obtaining these expressions. In order to use the recurrence relations for $V_k$, $A_k$, and $F_k$, the velocity and acceleration of the initial site, $V_0$ and $A_0$, must be provided, while the force associated with the site at the end of the final bond, $F_{n_r+1}$, is zero, since there is no joint associated with that site.

The purpose of the recurrence relations in Eqs. (18) and (24) is to provide expressions for $\{\dot{W}_k\}$, which, together with $A_0$, and assuming all the forces acting on the sites are known, can be integrated to solve for the chain dynamics; this is actually the opposite of the typical robotics problem, in which the goal is to determine the forces required to produce a particular robot arm trajectory.

### D. Stacked operators

Equations (17), (18), (24), and (25) can be rewritten in condensed, "stacked" form

$$V = \phi^T V + H^T \dot{W}, \quad (28)$$

$$A = \phi^T A + H^T \ddot{W} + X, \quad (29)$$

$$F = \phi F + M A + Y, \quad (30)$$

$$T = H F, \quad (31)$$

that combines the entire set of $k$ values. A quantity such as $V$ containing all the $V_k$ values for the chain is also referred to as a spatial vector, while, for example, the block matrix $\phi$ containing all the $\phi_{k,k+1}$ matrices is a spatial operator. The stacked formalism leads to a concise and elegant formulation of the problem, free from inundation by indices as is often the case in the robotics literature, e.g., Ref. [21].

The spatial operator approach was originally developed for the case of a fixed initial bond [16]—the base in the example of a robot arm—for which $V_0 = 0$, so that $\dot{W} = (\dot{\theta}_1, \ldots, \dot{\theta}_{n_r})^T$ is a vector with just $n_r$ components, and the other vectors and matrices are sized accordingly. In order to remove the fixed-base restriction [22], six extra DOFs are added to the problem by redefining $\dot{W} = (V_0, \dot{\theta}_1, \ldots, \dot{\theta}_{n_r})^T$ as a vector with $n_r + 6$ components; likewise for $\ddot{W}$. The size of the original $6n_r \times n_r$ block-diagonal matrix $H = \text{diag}(H_1, \ldots, H_{n_r})$ is increased to $6(n_r+1) \times (n_r+6)$ by including an extra $6 \times 6$ block $H_0 = I$, so that now $H = \text{diag}(I, H_1, \ldots, H_{n_r})$. The block-diagonal matrix $M$ is of size $6(n_r+1) \times 6(n_r+1)$; $\phi$ has the same size, and its only nonzero blocks are those to the immediate right of the diagonal. Namely $\{\phi_{01}, \ldots, \phi_{n_r-1,n_r}\}$. Vectors $V$, $A$, $F$, $X$, and $Y$, all have $6(n_r+1)$ components, e.g., $V = (V_0, \ldots, V_{n_r})^T$, and $T$ is organized in the same way as $\dot{W}$, with $n_r + 6$ components; $T_0 = 0$ because the special $k = 0$

joint exerts no torque. (Note that index order has been reversed from the original to make it more suitable for polymer use, and, for convenience, other aspects of the notation have been altered or simplified.)

The next step is to define the matrix

$$\Phi = (I - \phi)^{-1}, \quad (32)$$

which is also used in the alternative form, $\Phi = \Phi \phi + I$; because $\phi^{n_r+1} = 0$, Eq. (32) is equivalent to $\Phi = I + \phi + \phi^2 + \cdots + \phi^{n_r}$, which is an upper-triangular block matrix whose elements, each a $6 \times 6$ matrix, are

$$\Phi_{ij} = \begin{cases} I, & j = i, \\ \phi_{i,i+1}, & j = i+1, \\ \phi_{i,i+1} \cdots \phi_{j-1,j}, & j > i+1. \end{cases} \quad (33)$$

Then, in terms of $\Phi$, Eqs. (28)–(31) reduce to

$$V = \Phi^T H^T \dot{W}, \quad (34)$$

$$A = \Phi^T (H^T \ddot{W} + X), \quad (35)$$

$$T = \mathcal{M} \ddot{W} + H\Phi(M\Phi^T X + Y), \quad (36)$$

where

$$\mathcal{M} = H\Phi M \Phi^T H^T. \quad (37)$$

While $M$ is a sparse, $6(n_r+1) \times 6(n_r+1)$ block-diagonal matrix, $\mathcal{M}$ is only of size $(n_r+6) \times (n_r+6)$, but, although it is typically much smaller, it is fully populated. In principle, Eq. (36) can be numerically integrated to obtain $W$, and this is one of the approaches actually used in solving the problem, but the computational effort required for evaluating $\mathcal{M}^{-1}$ at each time step to obtain $\ddot{W}$ is of order $O((n_r+6)^3)$; for this reason such an approach is not practical for any but the shortest of chains. The alternative method, described below, requires a computational effort of order $O(n_r)$, together with what amounts to the inversion of a $6 \times 6$ matrix; clearly this will prove to be a far more efficient approach, even for relatively small $n_r$.

### E. Inversion of the mass matrix

As a preliminary step in obtaining an explicit expression for $\mathcal{M}^{-1}$ define [16] the $6 \times 6$ matrix $P_k$ in terms of $M_k$ as

$$P_k = \phi_{k,k+1}(I - G_{k+1}H_{k+1})P_{k+1}\phi_{k,k+1}^T + M_k. \quad (38)$$

In Eq. (38),

$$G_k = P_k H_k^T D_k^{-1}, \quad (39)$$

$$D_k = H_k P_k H_k^T, \quad (40)$$

where, for joints with a single DOF, $G_k$ is a six-component vector and $D_k$ is a nonzero scalar; note also that $P_k$ is symmetric. (The motivation for introducing $P_k$ is explained in

Ref. [16] and derives from the formal similarity of these equations with those used in the completely unrelated field of linear filtering.) Also define

$$\psi_{k,k+1} = \phi_{k,k+1}(I - G_{k+1}H_{k+1}) \qquad (41)$$

and substitute this in Eq. (38). The stacked versions of Eqs. (38)–(41) are

$$P = \psi P \phi^T + M, \qquad (42)$$

$$G = PH^T D^{-1}, \qquad (43)$$

$$D = HPH^T, \qquad (44)$$

$$\psi = \phi(I - GH). \qquad (45)$$

Matrices $P$ and $\psi$ are of size $6(n_r+1) \times 6(n_r+1)$, and $G$ is $(n_r+6) \times 6(n_r+1)$ and block-diagonal (thus the product $G_{k+1}H_{k+1}$ is square). Matrix $D$ is of size $(n_r+6) \times (n_r+6)$; its first $6 \times 6$ diagonal block corresponds to $D_0$, and the remaining $n_r$ diagonal elements are the scalars $D_k$. From Eqs. (42) and (45),

$$M = P - \phi P \phi^T + \phi GHP \phi^T, \qquad (46)$$

and so, by using Eq. (32),

$$\Phi M \Phi^T = P + \Phi \phi P + P \phi^T \Phi^T + \Phi \phi PH^T D^{-1} HP \phi^T \Phi^T. \qquad (47)$$

Substitute Eq. (47) in Eq. (37), then use $GD = PH^T$ from Eq. (43), together with Eq. (44), to obtain

$$\mathcal{M} = HPH^T + H\Phi \phi PH^T + HP \phi^T \Phi^T H^T$$
$$+ H\Phi \phi PH^T D^{-1} HP \phi^T \Phi^T H^T$$
$$= (I + H\Phi \phi G)D(I + H\Phi \phi G)^T. \qquad (48)$$

This alternative factorization of $\mathcal{M}$ is a product of three $(n_r+6) \times (n_r+6)$ matrices, unlike Eq. (37) that involves nonsquare matrices.

It is now a straightforward matter to invert $\mathcal{M}$. Use a special case of the Woodbury formula for the inverse of a matrix [23] $(I + Q_1 Q_2)^{-1} = I - Q_1(I + Q_2 Q_1)^{-1} Q_2$ to write

$$(I + H\Phi \phi G)^{-1} = I - H\Phi(I + \phi GH\Phi)^{-1} \phi G. \qquad (49)$$

By analogy with Eq. (32) for $\Phi$, define $\Psi = (I - \psi)^{-1}$; then from Eqs. (45) and (32),

$$\Psi^{-1} = \Phi^{-1} + \phi GH, \qquad (50)$$

so that $(I + H\Phi \phi G)^{-1} = I - H\Psi \phi G$. Thus the inverse of Eq. (48) is

$$\mathcal{M}^{-1} = (I - H\Psi \phi G)^T D^{-1}(I - H\Psi \phi G), \qquad (51)$$

and so, from Eq. (36),

$$\ddot{W} = (I - H\Psi \phi G)^T D^{-1}(I - H\Psi \phi G)[T - H\Phi(M\Phi^T X + Y)]$$
$$= (I - H\Psi \phi G)^T D^{-1}[T - H\Psi(\phi GT + M\Phi^T X + Y)], \qquad (52)$$

where Eq. (50) is used in simplifying $H(I - \Psi \phi GH)\Phi = H\Psi$. To eliminate $\Psi$, first rewrite Eq. (52) as

$$(I + H\Phi \phi G)^T \ddot{W} = D^{-1}[T - H\Psi(\phi GT + M\Phi^T X + Y)]. \qquad (53)$$

Next, use Eq. (42) with Eq. (32) to get

$$\Psi M \Phi^T = \Psi P(\phi^T \Phi^T + I) - \Psi \psi P \phi^T \Phi^T = \Psi P + P \phi^T \Phi^T. \qquad (54)$$

Then, using the transpose of Eq. (43), it follows that

$$(I + H\Phi \phi G)^T \ddot{W} = D^{-1}E - G^T \phi^T \Phi^T X, \qquad (55)$$

in which the forcelike quantities

$$E = T - HZ, \qquad (56)$$

$$Z = \Psi(\phi GT + PX + Y) \qquad (57)$$

have been defined. Rearranging Eq. (55) and using the expression for $A$ given in Eq. (35) leads to

$$\ddot{W} = D^{-1}E - G^T \phi^T \Phi^T(H^T \ddot{W} + X) = D^{-1}E - G^T \phi^T A. \qquad (58)$$

It is also possible to eliminate $\Psi$ from Eq. (57) by substituting $T$ from Eq. (56) to get $(I - \Psi \phi GH)Z = \Psi(\phi GE + PX + Y)$, and then using Eq. (50) to obtain

$$Z = \Phi(\phi GE + PX + Y). \qquad (59)$$

Explicit forms for the new recurrence relations embodied in Eqs. (58) and (59) are obtained by using Eq. (32) and reintroducing the $k$ indices,

$$Z_k = \phi_{k,k+1}(Z_{k+1} + G_{k+1}E_{k+1}) + P_k X_k + Y_k, \qquad (60)$$

$$\ddot{W}_k = D_k^{-1}E_k - G_k^T \phi_{k-1,k}^T A_{k-1}. \qquad (61)$$

These recurrence relations are used in opposite $k$ directions; they succeed in providing the required results without the need for explicit evaluation of the matrix inverse $\mathcal{M}^{-1}$ as implied by Eq. (36). It is for this reason that the method has not been referred to as an "inverse matrix method," a term sometimes seen in the literature, but rather a "rigid link" method, a far more apt descriptor.

The expressions given here describe the entire chain, but, provided the end joints are handled correctly, these results can be used for linear segments that form part of a larger assembly, allowing more complicated treelike structures to be treated. Furthermore, while the above formulation deals with the simplest case of a linear chain with a single torsional DOF per joint, it is readily extended to more complex joints, enabling, for example, the constant bond-angle condition to be eliminated by allowing two DOFs at each joint (an

alternative would be to decompose an individual joint into two coincident joints each with a single DOF).

## III. SIMULATION TECHNIQUES

### A. Linked-chain equations of motion

The recurrence relations used to propagate velocities, forces, and accelerations along the chain are as follows: The (translational and rotational) velocities $V_k$ are obtained by starting with $V_0$ and iterating Eq. (17),

$$V_k = \phi_{k-1,k}^T V_{k-1} + H_k^T \dot{W}_k, \quad k = 1, \ldots, n_r. \quad (62)$$

The forces (and torques), as represented by $E_k$, together with the matrices $D_k$ and $G_k$, are obtained by iterating Eqs. (38) and (60). For computational convenience, new quantities $A_k'$ and $Z_k'$ are introduced; then, starting with $P_{n_r+1} = 0$ and $Z_{n_r+1}' = 0$,

$$\left.\begin{aligned}
P_k &= \phi_{k,k+1}(I - G_{k+1}H_{k+1}) \\
&\quad \times P_{k+1}\phi_{k,k+1}^T + M_k, \\
D_k &= H_k P_k H_k^T, \\
G_k &= P_k H_k^T D_k^{-1}, \\
Z_k &= \phi_{k,k+1} Z_{k+1}' + P_k X_k + Y_k, \\
E_k &= T_k - H_k Z_k, \\
Z_k' &= Z_k + G_k E_k,
\end{aligned}\right\} k = n_r, \ldots, 0. \quad (63)$$

Finally, the values of $\ddot{W}_k$ (or $\ddot{\theta}_k$) are determined by starting with $A_0$ (its evaluation is discussed below), and iterating Eqs. (18) and (61),

$$\left.\begin{aligned}
A_k' &= \phi_{k-1,k}^T A_{k-1} \\
\ddot{W}_k &= D_k^{-1} E_k - G_k^T A_k' \\
A_k &= A_k' + H_k^T \ddot{W}_k + X_k
\end{aligned}\right\} k = 1, \ldots, n_r. \quad (64)$$

These recurrence relations, which are readily transformed into a suitable computer program, imply a series of operations (multiplications and additions) involving $6 \times 6$ matrices and six-component vectors, but the total computational effort is only of order $O(n_r)$.

Recall that the $k=0$ joint has six DOFs, and also that $H_0 = I$, $X_0 = 0$, and $\ddot{W}_0 = A_0$. Now, because $A_{-1} = 0$, it follows from Eq. (64) that $A_0 = D_0^{-1} E_0$, and since $T_0 = 0$,

$$D_0 A_0 = -Z_0, \quad (65)$$

where both $D_0$ and $Z_0$ have already been determined (above). Thus $A_0$ can be evaluated numerically by solving the set of six linear equations contained in Eq. (65) using the standard LU decomposition method [23]; the computational effort required for this initial joint is fixed and independent of $n_r$.

### B. Leapfrog integration and rigid-body equations

The familiar leapfrog method for integrating the MD translational equations of motion—which is algebraically equivalent to the Verlet method [24]—is usually expressed in a form where the coordinates and velocities are evaluated at alternate half time steps [15]. This minor inconvenience can be avoided by using a slightly modified form that breaks the integration procedure for a single time step into two parts. Prior to computing the latest acceleration (**a**) values, update the velocities (**v**) by a half time step using the previous accelerations, and then update the coordinates (**r**) by a full time step using these intermediate velocity values,

$$\mathbf{v}(t+h/2) = \mathbf{v}(t) + (h/2)\mathbf{a}(t), \quad (66)$$

$$\mathbf{r}(t+h) = \mathbf{r}(t) + h\mathbf{v}(t+h/2). \quad (67)$$

In the case of the polymer chain, this procedure is applied to the translation coordinates of the $k=0$ site and (in scalar form) to each of the dihedral angles $\theta_k$; the treatment of the angular coordinates associated with the $k=0$ site, below, employs a related approach for dealing with the rotational equations. Next, use the new coordinates (and velocities if needed) to compute the latest acceleration values, then update the velocities over the second half time step,

$$\mathbf{v}(t+h) = \mathbf{v}(t+h/2) + (h/2)\mathbf{a}(t+h). \quad (68)$$

In the linked-chain formulation, the initial bond of the chain is treated as a rigid body; the influence of the rest of the chain on it has already been taken into account and is contained in the force and torque transmitted through the first internal joint. There are a number of ways of describing the orientation of a rigid body [10]: Euler angles have proved very useful for analytic purposes because of their intuitive nature, but owing to a potentially singular matrix that appears in the equations of motion they are not the preferred method for dealing with numerical problems. Quaternions have achieved popularity because of their singularity-free nature, but their normalization must be preserved against a small but persistent numerical drift [25,15]. A more recently proposed alternative is to regard the complete rotation matrix as the dynamical variable; this is the representation that will be used here, since the integration scheme [19]—which is based on operator splitting and maintains time reversibility—is just another instance of the leapfrog method.

In the original description [19], vector components were expressed in the principal-axis frame of the body. Since the chain dynamical problem as a whole is solved in the space-fixed frame, the corresponding form of the rotational equations will be described here. If **R** denotes the rotation matrix of a rigid body, then the first part of the leapfrog integration step consists of a half-time-step update of the angular velocities,

$$\boldsymbol{\omega}(t+h/2) = \boldsymbol{\omega}(t) + (h/2)\boldsymbol{\alpha}(t), \quad (69)$$

where $\boldsymbol{\alpha} \equiv \dot{\boldsymbol{\omega}}$, followed by a full-time-step update of **R** using a symmetric product of matrices describing a series of small partial rotations,

$$R^T(t+h) = U_1 U_2 U_3 U_2 U_1 R^T(t), \qquad (70)$$

where, for convenience, the transpose of $R$ is treated. Note that for the linked chain, the rigid body is associated with the $k=0$ joint, so that $R \equiv R_0$. Each of the matrices

$$U_1 = U_x(\omega_x h/2), \quad U_2 = U_y(\omega_y h/2), \quad U_3 = U_z(\omega_z h) \qquad (71)$$

describes a rotation about a single axis and is evaluated in the space-fixed frame. For small angles, they can be approximated in a way that preserves orthogonality, e.g.,

$$
U_x(\theta) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos\theta & -\sin\theta \\ 0 & \sin\theta & \cos\theta \end{pmatrix}
$$

$$
\approx \begin{pmatrix} 1 & 0 & 0 \\ 0 & \dfrac{1-\theta^2/4}{1+\theta^2/4} & \dfrac{-\theta}{1+\theta^2/4} \\ 0 & \dfrac{\theta}{1+\theta^2/4} & \dfrac{1-\theta^2/4}{1+\theta^2/4} \end{pmatrix}. \qquad (72)
$$

The second part of the leapfrog step is

$$\omega(t+h) = \omega(t+h/2) + (h/2)\alpha(t+h). \qquad (73)$$

In the case of a single rigid body, the angular acceleration is determined from the torque $\tau$, namely, $\alpha(t+h) = \mathcal{I}^{-1}\tau(t+h)$, whereas for the linked chain this treatment is only required for the $k=0$ joint, and $\alpha$ is obtained by solving Eq. (65); the reason rigid bodies are usually treated in the body-fixed principal-axes frame is to ensure the diagonality of $\mathcal{I}$, a consideration that is not relevant here.

The complete procedure for a single time step can be summarized as the following sequence of operations. Integrate (first part) to obtain base velocities and coordinates, and joint angular velocities and angles; determine site velocities, Eq. (62); evaluate site coordinates, Eqs. (1)–(3); compute external forces and torques, and other necessary quantities; determine joint forces, Eq. (63); solve Eq. (65) for the base acceleration; determine joint accelerations, Eq. (64); integrate (second part) to obtain base velocities and joint angular velocities.

### C. Polymer chain model

Two kinds of interactions are required in this model—excluded volume and torsion. The former is provided by a pair interaction that prevents overlap of the atoms (or atom groups) located at the chain sites. Here a simple soft-sphere repulsion, based on the Lennard-Jones potential with a short-range cutoff, is all that is required: for a pair of atoms located at $r_i$ and $r_j$, where $r_{ij} = r_i - r_j$, and $r_{ij} = |r_{ij}|$, the potential is

$$
u_{ss}(r_{ij}) = \begin{cases} 4\epsilon[(r_{ij}/\sigma)^{12} - (r_{ij}/\sigma)^6], & r_{ij} < r_c, \\ 0, & r_{ij} \geq r_c, \end{cases} \qquad (74)
$$

with a cutoff $r_c = 2^{1/6}\sigma$ (nearby pairs of atoms that are prevented from approaching too closely because of geometrical restrictions need not be considered). Should a pairwise attraction between particular pairs of distant chain atoms be required (as will be the case later on), it can be obtained from Eq. (74) by simply increasing $r_c$. The pair forces derived from this potential, and their associated torques, contribute to $f_k^e$ and $n_k^e$ in Eqs. (9) and (10).

The torsional potential associated with the dihedral angles $\theta_k$ has the simple form

$$u_t(\theta_k) = -u_k \cos(\theta_k - \theta_k^{(0)}), \qquad (75)$$

where $\theta_k^{(0)}$ is the dihedral angle that produces a ground state having the correct helical twist, and $u_k$ is the interaction strength. The torque appearing in Eq. (13) is

$$t_k = u_k \sin(\theta_k - \theta_k^{(0)}), \qquad (76)$$

a result whose simplicity stands in sharp contrast to the intricate vector algebra associated with torque calculations when working in Cartesian coordinates [15].

For the chains considered here it is assumed all $|b_k| = b$, $\alpha_k = \alpha$, $\theta_k^{(0)} = \theta^{(0)}$, and, except for the later twin-helix studies where selected $u_k = 0$, all $u_k = u^{(0)}$. Since the torsion also acts at the first internal joint, it is necessary to add an extra site and bond to the chain (effectively with an index "$-1$") to make this torsion term meaningful; the first three sites of the modified chain form a rigid unit (the extra bond does not alter the preceding analysis) and the chain length is increased by unity.

A spherical mass element (with a finite moment of inertia about its own center of mass) is associated with each site; for bonds with $k > 0$, the mass is attached to the far ($k+1$ site) of the bond, while the $k=0$ bond, as explained above, has three masses associated with it. The components of the inertia tensor in Eqs. (26) and (27) are

$$
(\mathcal{I}_k)_{ij} = \begin{cases} \sum_{\kappa \in k} m_\kappa (r_\kappa^2 - r_{\kappa i}^2), & i = j, \\ -\sum_{\kappa \in k} m_\kappa r_{\kappa i} r_{\kappa j}, & i \neq j, \end{cases} \qquad (77)
$$

where the sum (or volume integral) is over all mass elements $\kappa$ fixed to bond $k$, and coordinates are relative to the center of mass of each bond in the space-fixed frame.

### D. Order parameter

While the appearance of an ordered helical structure, even one with the occasional defect, is easily recognized visually, in order to facilitate statistical analysis of the behavior it is important to be able to quantify the degree of order present in the chain. Let $d_k = b_{k-1} \times b_k$, then

$$S = \frac{1}{n_r} \left| \sum_{k=1}^{n_r} \hat{\boldsymbol{d}}_k \right| \qquad (78)$$

defines an order parameter that measures the long-range order present in the folded structure based on the orientation of the helical turns; for a single, well-formed helix, $S$ should have a value close to unity. A slightly modified version of $S$ will be introduced later for studying twin-helix structures.

This definition of $S$ is particularly useful for detecting structures consisting of two or more helical domains with axes aligned in different directions due to a localized defect of the type seen in helically wound telephone and electrical cords. Since the correct helicity (or "handedness") is built into the interactions, it is unlikely that segments of opposite helicity will independently nucleate at separate locations but, as the chain collapses, individual turns with the wrong twist can become trapped in the structure. These defects are capable of traveling along the chain, but this is a slow process, and the direction of motion is random unless close to the chain end. There are instances where the definition of $S$ in Eq. (78) can give an incomplete picture; if a wrong turn occurs very close to the chain end, its effect on $S$ will be minimal, and even a perfectly formed helix is subject to low frequency bending motion. Other order parameters can be defined that are of a more short-range nature; for example, a simple count of the number of pairs of chain sites lying within a specified range (i.e., the number of "contacts" between adjacent turns of the helix) divided by the maximum possible value, but for long chains the tolerance in the threshold required to accommodate thermal fluctuations might allow significant changes in the helical-axis direction to go undetected.

## IV. RESULTS

### A. Simulation details

One of the more prominently recognizable structural motifs found in proteins is the ($\alpha$) helix. The helix, because of its uniformity along the longitudinal axis, is a particularly simple structure to specify, and both Monte Carlo and MD helix-folding simulations based on the complex potentials designed for protein modeling have been carried out, e.g., Refs. [18,26]. Complex potentials have also been used in MD studies of reversible folding processes that involve helical states [27]. Since the complexity of these potentials is not obviously essential for a basic understanding of generic folding phenomena, the present simulations are based on the much simpler model and potentials described previously. Indeed, an analogous approach has been employed experimentally [28] in a study of helix formation in synthesized non-biological chain molecules, where the interactions are simpler than in proteins (in particular, there are no hydrogen bonds).

The importance of examining simple structures, such as the helix, is that the process by which ordered arrangements emerge from randomly coiled states is likely to capture something of the essence of real protein folding, such as the cooperativity of the folding process, the role of nucleation sites, the degree to which folding is able to proceed to completion, and the steric and topological effects of excluded volume. The key question, of course, is whether the intrinsic time scales of these processes are sufficiently small for simulation to be computationally feasible, an issue addressed by the results presented here. While the present model is admittedly a mere caricature of the detailed models normally employed in studies of individual proteins, it has two undeniable advantages, namely, a known native ground state compatible with the interactions, and sufficiently modest computational requirements that MD simulation is able to encompass the time interval required for major conformational change. More complex protein structures also display certain common characteristics, and ought to be accessible to simulations of this type; it is, however, essential to eliminate any ambiguity from the ground state, something that nature itself has presumably achieved in the interests of efficiency and reliability.

Each simulation run considers a single chain constructed as described earlier. The absence of a solvent, apart from changing the time scales, should not alter the outcome; indeed many, if not most, protein simulations avoid introducing an explicit solvent for reasons of computational efficiency. The simulation is begun at a relatively high temperature, so that the kinetic energy is sufficiently large to surmount the torsional potential barriers. The initial chain configuration is a large loop extending across the simulation cell, with a very slight helicity to prevent any overlap; initial dihedral angles are chosen so that locally, the conformation is almost a planar zigzag state. The joint angular velocities are assigned random values corresponding to the starting temperature, and memory of this initial state rapidly vanishes early in the simulation. The temperature is gradually reduced by a factor slightly less than unity at regular intervals until, towards the end of the run, very little kinetic energy remains in the system. The simulation region is bounded by hard, reflecting walls; while there are occasional wall collisions, this has little influence on the overall behavior. (The alternative would be to use periodic boundaries, which for a simulation cell not large enough to contain the chain in a fully stretched state, would be subject to chain wraparound effects; while these are also unlikely to affect the overall behavior, they can prove visually confusing given the importance of computer-generated visualization in this work.)

The gradual cooling that is imposed throughout the run plays several distinct roles. During the early stage it is used to drive the chain from a totally random state to one in which the torsional potential begins to have some influence over the dihedral angles. Then, as the temperature is reduced further, an increasing degree of local order emerges and precursors to long-range order appear, either as a consequence of the merging of separate ordered domains, or the spread of order from a nucleation region (or a combination of both processes); during this stage the imposed cooling performs a task normally the responsibility of the solvent, namely, the removal of excess potential energy as the chain evolves towards states of lower energy. Once the chain has reached a state consisting mainly of helical segments, possibly separated by small misfolded regions that have become trapped, the purpose of further temperature reduction is to gradually freeze out thermal fluctuations—without further major struc-

tural change—in order to allow evaluation of the long-range order parameter $S$ (the measure of success of the folding process); the latter part of this cooling stage is not intended to imitate any real physical process.

The simulations use standard, reduced MD units, in which all distances and energies are expressed in terms of the Lennard-Jones parameters $\sigma$ and $\epsilon$, respectively; mass is expressed in terms of the monomer mass $m$ and, consequently, the unit of time is $\sqrt{m\sigma^2/\epsilon}$. Temperature and energy are made numerically identical by setting the Boltzmann constant $k_B$ to unity. In terms of these units the parameters used in the runs are as follows: The bond length $b=1.3$, a value sufficiently short to prevent the chain crossing itself, the bond angle $\alpha$ and the preferred dihedral angle $\theta^{(0)}$ are chosen to produce helices with periodicity six, and the torsional potential strength $u^{(0)}=5$. In the studies of twin helices, the cutoff in the attractive interaction, based on Eq. (74), occurs at $r_c=2.2$. The initial temperature is 4 (corresponding to a kinetic energy per DOF of 2) and the final temperature is $10^{-3}$; temperature is reduced by rescaling all velocities and angular velocities by a factor $f_T$ every 4000 time steps, with $f_T=0.95$ or 0.97. The runs reported here are each of length $4$–$8\times10^5$ steps; the integration time step (in MD units) is $h=4\times10^{-3}$. In order to produce reliable statistics, a large number of runs were carried out for each case studied; the runs differed in the choice of initial random values for $\{\dot{\theta}_k\}$.

### B. Folding to a single helix

Measurements were made of the long-range order parameter $S$ and the total energy, the latter a sum over contributions from the soft-sphere pair interactions, Eq. (74), the torsional terms, Eq. (75), and the kinetic energy. The measurements involved 400 independent runs for each of several chain lengths $L$ and different cooling rates. These quantitative results were complemented by an interactive graphical version of the simulation program that provided real-time visual monitoring of the folding process; in addition to learning about any potential obstructions to complete folding, the ability to observe chains directly also helped when choosing a cooling rate sufficiently fast for folding to proceed to completion, but not too fast for an excessive number of defects to become trapped in the nascent structures.

The viability of the underlying approach depends on whether it can actually produce correctly structured helices. The first series of results measures the fraction of chains that successfully fold into a helical state, and the manner in which the success rate depends on $L$ and the cooling rate. A summary appears in Table I; $L$ ranges from 18 to 90, which, since the helix period is six, corresponds to 3–15 full helical turns.

Owing to the large number of runs it is not possible to provide a detailed history of each, so a quantitative measure of folding success must be introduced. A successfully folded helix is deemed to be one for which $S>0.88$ at least once during the last $1.2\times10^5$ steps of the run (measurements are made every 4000 steps); by this stage of the run the system has reached a comparatively low temperature, so that further substantial conformational changes are unlikely. Visual

TABLE I. Details of helix folding runs discussed in the text.

| Length ($L$) | Turns | $f_T$ [a] | Steps ($\times10^3$) | Success [b] |
|---|---|---|---|---|
| 18 | 3 | 0.95 | 400 | 1.00 |
| 36 | 6 | 0.95 | 400 | 0.94 |
| 54 | 9 | 0.95 | 400 | 0.85 |
| 54 | 9 | 0.97 | 800 | 0.94 |
| 72 | 12 | 0.95 | 400 | 0.69 |
| 72 | 12 | 0.97 | 800 | 0.91 |
| 90 | 15 | 0.97 | 800 | 0.85 |

[a]Cooling factor.
[b]Criterion for successful helix formation is defined in the text.

analysis confirms that, for the cases considered, this threshold for $S$ provides a quite reliable estimator; it tends to be sensitive to defects in the helical structure, while allowing for the fact that a properly folded helix may still have some residual curvature along its major axis.

It is clear from Table I that a high success rate for helix production is achieved. Two trends are apparent in the results, neither of them unexpected. For a given $f_T$, longer chains are less likely to fold properly than shorter chains, and, for a given $L$, a larger $f_T$ (corresponding to slower cooling) raises the success rate. Thus the longer the chain, the slower the desired cooling rate; additional runs with faster cooling confirm this observation. The longest of the chains folds to a helix with 15 turns, which, considering the potential for defects, represents a significant victory of energy over entropy.

The rate at which chains approach the helically ordered state can be studied by monitoring the mean values of $S$, as well as the negative of the total energy (which is dominated by the torsional component when in the folded state); these quantities provide measures of the long- and short-range order, respectively. The results, normalized per DOF, for 54- and 90-site chains, averaged over all 400 runs, are shown in Figs. 3 and 4. The overall results are divided into two groups, depending on whether the chain is classified as having folded
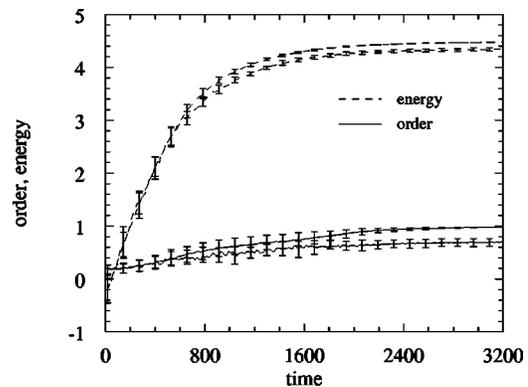


FIG. 3. Averaged order parameter and (negative) total energy per DOF as functions of time (in dimensionless MD units) for chains with $L=54$; the contributions of chains that do and do not fold correctly appear in separate curves, with the upper curve in each case corresponding to the successful folders.
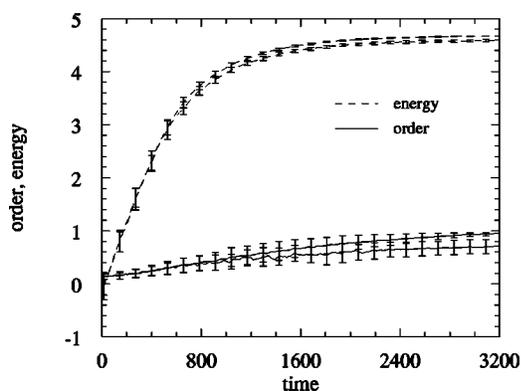
FIG. 4. Order parameter and energy for $L=90$ (similar to Fig. 3).

successfully or not, and error bars indicate the standard deviations of the measurements. In each case it is the upper curve that represents the average for the successfully folded chains, and it has the smaller error bars.

An alternative estimate of the rate at which folding proceeds is based on the time dependence of the fraction of chains in a helical state, relative to all those that eventually succeed in reaching this state. This provides information about when, assuming a chain folds successfully, the appearance of helical order actually occurs. Figure 5 shows these cumulative distributions for the different $L$, each at the slowest cooling rate considered. The cooling rate clearly affects the results, as can be seen from the separation of the two groups of curves that are based on different rates (see Table I); for a given cooling rate, the folding speed tends to drop as the chains become longer (the slight crossover of the $L=72$ and 90 curves is probably not significant).

A more detailed examination of final-state conformations is based on histograms of the $S$ distribution, using measurements made over the last $1.2 \times 10^5$ steps. These results appear in Fig. 6 for several $L$ values (at the slowest cooling rate). There are two separate curves for each $L$, one showing the spread of $S$ for those chains that satisfied the folding criterion at least once during this measurement period, and a broader, much lower curve for the chains that did not. The former set of distributions become broader with increasing $L$;
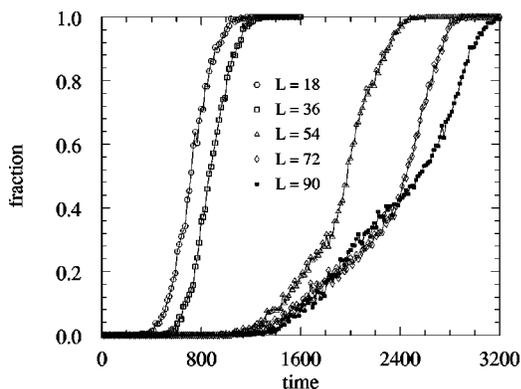
FIG. 5. Cumulative distributions of chains in the folded state as a function of time, for different lengths ($L$).
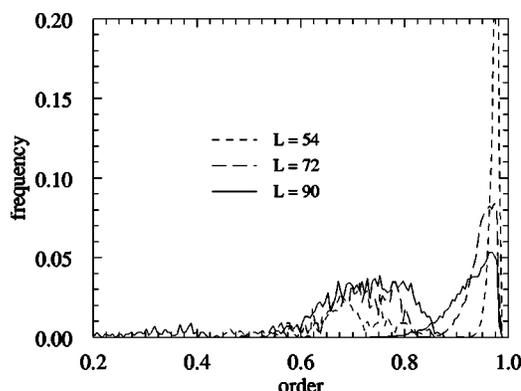
FIG. 6. Order parameter distributions; separate curves show results for chains that did (peaks on the right) and did not (multiplied by a factor of 10) fold correctly.

there are several contributing causes for this, including slower folding rates (all the runs were of equal length), chains not managing to fold successfully but having one or more intermediate $S$ values which passed the test, and the increasing effect of bending along the helical axis. The latter, broader distributions (scaled up by a factor of 10 to make the details visible) are due both to the many defective structures possible, each with its own spread of $S$ values, and also to the defects themselves reducing the structural rigidity and so raising the susceptibility to slow thermal vibration. Only for the longest chains is there any overlap of the curves, and even then it is minimal.

The best way to follow the folding process is by viewing animated sequences of images taken at various points during the run; some sequences can actually be generated while running the simulation interactively, if the computations proceed sufficiently rapidly. Here, due to the limitations of the printed page, a selection of static images must suffice. One could attempt a verbal description of what transpires but, as was the case in Ref. [18], there are no obvious features shared by the individual folding trajectories. Even if certain common characteristics do exist, the strong random conformational fluctuations make their observation difficult; a systematic, quantitative means for identifying pathways, that extends ideas used for equilibrium states [9], might prove helpful in this task.

Figure 1, which appeared early in the paper, shows an image of a typical, well formed, almost straight helix obtained in one of the runs, while Fig. 2 shows a random chain configuration observed near the start of a run, both for $L=90$ chains. For clarity, these and subsequent pictures represent the chains by their tubular envelopes, rather than by showing individual, partially overlapped spheres represent-
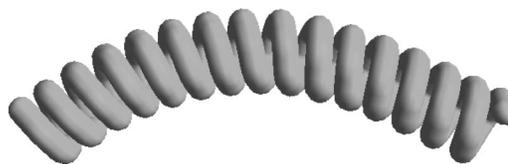
FIG. 7. Correctly folded helix ($L=90$) with residual curvature; this is the most extreme case of bending observed.
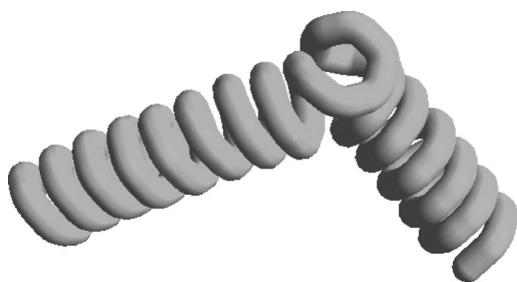
FIG. 8. Incorrectly folded state with a single defect.



FIG. 10. Folded state with a localized intertwined defect.

ing the monomers. The tube thickness corresponds to unit diameter, slightly less than the soft-sphere interaction cutoff to ensure that a small amount of space remains visible between adjacent turns of the helix. The maximum amount of residual curvature that was observed in the backbones of properly folded helices is apparent from Fig. 7; this example actually meets the criterion for folding success (as indeed it should).

While the majority of runs (85% for chains with $L=90$, and an even higher proportion for smaller $L$, see Table I) result in a correctly folded helix, examination of the kinds of defects that appear in the final states of those runs that fail to fold properly is an informative exercise. The first such picture, Fig. 8, is of an $L=90$ chain with two helical regions separated by a single defect. The defect is essentially a single loop of the helix with a reverse fold that became frozen in place during the cooling process. This is the most frequent type of defect, and its location can be anywhere in the chain, even right at the end.

Less frequent are chains with two spatially separated defects, as shown in Fig. 9. More extreme, but very rare examples of other kinds of defects appear in Figs. 10 and 11. These show what can happen when the chain starts to become entangled with itself; in the first case the problem is localized, but in the second (the only example of its kind observed) there is a relatively large loop trapped by the entanglement. The fact that these defects are relatively infrequent, and that even this low level of failure can be reduced by lowering the cooling rate still further, attests to the robustness and reliability of the folding process.

### C. Folding to a pair of helices

The helix formation described above is obtained in a study of homopolymers where, due to the uniformity of the chain, the only kind of repeating structure that can be produced is the helix (the planar zigzag conformation is a degenerate case). In a protein context, this kind of structure is
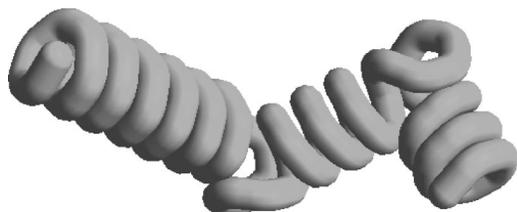
classed as secondary because helices often serve as structural components in more complex assemblies. Globular proteins are characterized by at least one additional level in the structural hierarchy, namely, tertiary structure. A model capable of demonstrating tertiary structure requires at least some differentiation among the chain sites that breaks the translational invariance. Building on the helix-forming model investigated here, the next stage of complexity is a packed assembly of helices, a structure that incorporates both the secondary and tertiary levels of the hierarchy. The simplest way to design such a structure is to include nonlocal, attractive forces between selected pairs of chain sites; here, the pairs involved are located at chain positions that will be brought into proximity following the collapse into a state with two adjacent helices aligned in antiparallel directions. Such highly specific interactions are reminiscent of an approach used for lattice protein models [29]. The overall simplicity should be contrasted with the highly detailed model, complete with solvent, used in an MD study of the *unfolding* of a three-helix bundle [30]; the folding of such bundles has been studied [31] using a continuum version of the simplified specific-interaction approach [29] by means of discrete-event MD [32,15].

Choosing the interactions to produce a twin-helix structure is accomplished as follows. For a homogeneous chain with $L=n_b+1$ sites, in which the periodicity of the helix is $p$, the ground state consists of $n_t=L/p$ turns. Now assume that $n_t$ is an odd number and choose the interactions appropriate for a pair of adjacent helices, each with $(n_t-1)/2$ turns, joined by a "bridging" chain segment of length $p$. All that remains is to identify the pairs of sites in the two helical regions that must attract; these are just neighboring sites in adjacent turns of one of the helical segments, matched with



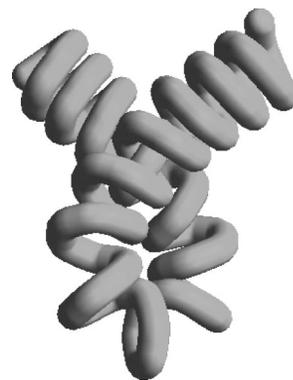FIG. 9. Incorrectly folded state with two separate defects.



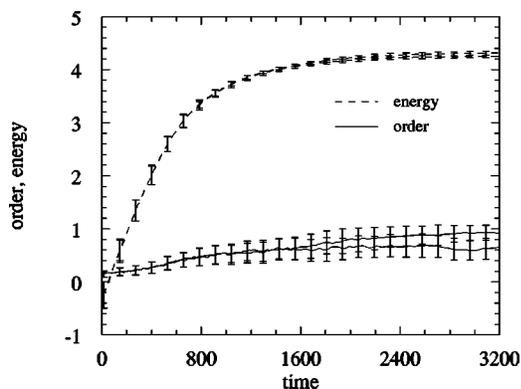FIG. 11. Folded state with a complex defect involving a large loop.

FIG. 12. Averaged order parameter and energy as functions of time for chains that form antiparallel helix pairs ($L=78$).

the corresponding sites, in reverse order, of the other. The strength of the attractive potential responsible for the tertiary structure, which is based on Eq. (74), is $0.2u^{(0)}$; it is weaker than the torsion, but the question of whether tertiary structure formation is the beneficiary or the cause of secondary structure formation [2] is not addressed here. In these exploratory computations, the torsional interactions along the bonds in the bridging segment are set to zero for simplicity; such interactions could actually be used to assist the folding and will be the subject of future study.

The definition of the long-range order parameter $S$, Eq. (78), must be modified to reflect the structure of the anticipated collapsed state. The partial contributions to $S$ of the two helical segments are now combined with opposite signs, and contributions from the bridging region ignored; this provides a reasonably sensitive, but unambiguous, measure of folding success. Figure 12 shows how both the modified $S$ and the energy vary with time, for $L=78$ chains (corresponding to a folded state consisting of a pair of six-turn helices), using runs whose details are otherwise similar to those for $L=90$. Based on visual analysis of the behavior, a different definition of what constitutes successful folding is needed here, namely, that the value of the modified $S$ must now exceed 0.95 for folding to be considered successful; using this criterion the success fraction was found to be 0.66.

Figure 13 shows an example of a successfully folded helix pair, an ordered conformation in which both secondary
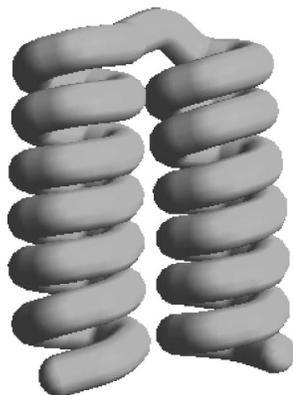


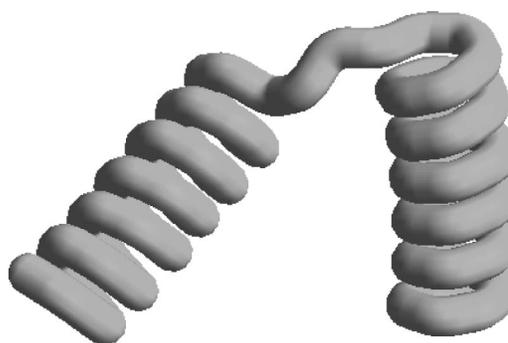FIG. 13. Well-formed pair of helices with antiparallel alignment.



FIG. 14. Pair of helices that have failed to align.

and tertiary structural elements are manifest; two thirds of the runs ended in this state. The failure to fold properly was generally not due to defects in the individual helical segments, but because the two secondary components did not succeed in aligning correctly; an example of such an outcome is shown in Fig. 14. The attractive forces become much more effective once the helical segments have formed; this is due to the linear arrangement of the attraction sites enabling them to function cooperatively, an effect that may be reflected in real proteins with prominent secondary-structural features. As a result of the nature of the interactions and the relative interaction strengths, the helical segments form first, essentially unimpeded, and only then do they attempt to align. The failure to align here is a symptom of the absence of any driving force for bringing the helices together; there is no torsional preference in the bridging segment and the range of the interhelix attraction is too short to be felt if the helical segments are well separated. Changes to either or both these aspects of the potential should alter the behavior, but care is required to avoid hindering the helix formation process in any way.

## V. CONCLUSIONS

The present paper has focused on both methodology and results. A formalism developed for the dynamics of robotic manipulators and other coupled mechanical systems—that provides a convenient and direct representation of the dynamics of bodies connected by rigid links with restricted degrees of freedom—has been utilized in a polymer context, with the clear implication that existing methods based on geometric constraints may be redundant in many instances. Since the treatment also involves dealing with rigid-body dynamics, a computationally more effective method than the often-used quaternion approach is also employed.

The results of an extensive series of MD simulations demonstrate that homopolymer chains with suitable torsional interactions consistently collapse into well-formed helices; the probability of localized defects being frozen into the structure depends on the cooling rate, and it can be reduced to a very low level by cooling sufficiently slowly. In order to demonstrate that the present simplified approach is relevant to protein folding, heterogeneous chains, with interactions favoring the development of antiparallel pairs of helices, were shown to produce coexisting secondary and tertiary structural features.

The order parameters introduced to quantify the degree of folding were tailored to capture the structural order present in the final state of the polymer. To study the details of folding pathways, other order parameters (or reaction coordinates) that capture features present in the intermediate states, but not necessarily in the final state, could be defined. In the twin-helix case, for example, a simple sum of the absolute values of $S$, evaluated separately for each helix-forming segment of the chain, might prove useful, since this quantity reaches its maximum upon completion of secondary structure formation, and is not seriously affected by subsequent rearrangement at the tertiary level.

The apparent success of the MD approach to chain folding used here is important for another reason. The widely cited Levinthal "paradox" [5] implies that since the number of states accessible to a protein grows exponentially with residue count, the time required for even a small protein to seek out its native state is, for practical purposes, infinite. Since nature does not suffer from this problem, the implication is that substantial portions of the folding process occur along certain well-characterized pathways; thus the molecules do not really wander almost aimlessly through conformation space, and hence there is no paradox. In order to begin to simulate such processes it is necessary to resort to a computationally efficient model, with realistic dynamics and a unique but readily determined low-energy "native" state; this is precisely what has been accomplished in the present work.

The type of model introduced here provides a starting point for exploration in several directions. While the interactions were weighted to construct the secondary helix structure prior to forming features at the tertiary level, a change in the relative strength of the interactions would allow aspects of both levels of organization to appear concurrently. Chains could be designed to fold into other idealized compact structures, such as the packed cube used in some lattice studies [2,3], or sheetlike conformations that also represent important secondary structure components; furthermore, packed states with different degrees of accessibility could provide useful information on how this feature influences folding success. The interactions can be modified and new types of interactions added; polymer topology can be changed by the addition of side chains corresponding to residues with extended structure. Common to all these enhancements is that the model must always be designed with a known lowest-energy state, and in this respect the approach differs from many other types of protein simulation. While models of this kind are perhaps limited in the kinds of questions they can address, there are more than enough issues requiring attention where they can prove helpful.

In a sense, the role played by such highly simplified models is analogous to the Ising model of ferromagnetism [33]; while it is not usually claimed that an Ising spin system accurately represents a real magnetic material (or, for that matter, any other kind of real physical system, when used for other kinds of problems such as lattice gases) it does, however, capture a great deal of the essence of the problem, to an extent that the study of Ising and related models has resulted in important advances, both for spin systems in particular, and for statistical mechanics and critical phenomena in general. Proteins can also be modeled with a high level of detail and specificity, but the tradeoff is that only short trajectory segments can be followed with an investment of a reasonable amount of computing effort; hopefully, extensive studies of simplified polymer models of the kind examined here, in which the design is tailored to reproduce certain generic aspects of macromolecular behavior, will achieve greater popularity as their usefulness becomes established.

[1] C. L. Brooks III, M. Karplus, and B. M. Pettit, *Proteins: A Theoretical Perspective of Dynamics, Structure, and Thermodynamics* (Wiley, New York, 1988).

[2] K. A. Dill, S. Bromberg, K. Yue, K. M. Fiebig, D. P. Yee, P. D. Thomas, and H. S. Chan, Protein Sci. **4**, 561 (1995).

[3] E. I. Shakhnovich, Curr. Opin. Struct. Biol. **7**, 29 (1997).

[4] C. L. Brooks III, Curr. Opin. Struct. Biol. **8**, 222 (1998).

[5] V. S. Pande, A. Y. Grosberg, and T. Tanaka, Rev. Mod. Phys. **72**, 259 (2000).

[6] Y. Duan and P. A. Kollman, Science **282**, 740 (1998).

[7] H. J. C. Berendsen, Science **282**, 642 (1998).

[8] M. Levitt and A. Warshel, Nature (London) **253**, 694 (1975).

[9] G. M. Crippen and Y. Z. Ohkubo, Proteins **32**, 425 (1998).

[10] H. Goldstein, *Classical Mechanics*, 2nd ed. (Addison-Wesley, Reading, MA, 1980).

[11] J.-P. Ryckaert, G. Ciccotti, and H. J. C. Berendsen, J. Comput. Phys. **23**, 327 (1977).

[12] G. Ciccotti, M. Ferrario, and J.-P. Ryckaert, Mol. Phys. **47**, 1253 (1982).

[13] H. C. Anderson, J. Comput. Phys. **52**, 24 (1983).

[14] R. Edberg, D. J. Evans, and G. P. Morriss, J. Comput. Phys. **84**, 6933 (1986).

[15] D. C. Rapaport, *The Art of Molecular Dynamics Simulation* (Cambridge University Press, Cambridge, 1995).

[16] G. Rodriguez and K. Kreutz-Delgado, IEEE Trans. Rob. Autom. **8**, 65 (1992).

[17] A. Jain, J. Guid. Control Dyn. **14**, 531 (1991).

[18] R. A. Bertsch, N. Vaidehi, S. L. Chan, and W. A. Goddard III, Proteins **33**, 343 (1998).

[19] A. Dullweber, B. Leimkuhler, and R. McLachlan, J. Comput. Phys. **107**, 5840 (1997).

[20] A. Jain, N. Vaidehi, and G. Rodriguez, J. Comput. Phys. **106**, 258 (1993).

[21] K. S. Fu, R. C. Gonzalez, and C. C. G. Lee, *Robotics: Control, Sensing, Vision, and Intelligence* (McGraw-Hill, New York, 1987).

[22] A. Jain and G. Rodriguez, IEEE Trans. Rob. Autom. **11**, 585 (1995).

[23] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. R. Flannery, *Numerical Recipes in C: The Art of Scientific Computing*, 2nd ed. (Cambridge University Press, Cambridge, 1992).

[24] L. Verlet, Phys. Rev. **159**, 98 (1967).

[25] D. J. Evans, Mol. Phys. **34**, 317 (1977).

[26] D. C. Rapaport and H. A. Scheraga, Macromolecules **14**, 1238 (1981).

[27] B. L. de Groot, X. Daura, A. E. Mark, and H. Grubmüller, J. Mol. Biol. **309**, 299 (2001).

[28] J. C. Nelson, J. G. Saven, J. S. Moore, and P. G. Wolynes, Science **277**, 1793 (1997).

[29] Y. Ueda, H. Taketomi, and N. Go, Biopolymers **17**, 1531 (1978).

[30] E. M. Boczko and C. L. Brooks III, Science **269**, 393 (1995).

[31] Y. Zhou and M. Karplus, Nature (London) **401**, 400 (1999).

[32] Y. Zhou, M. Karplus, J. M. Wichert, and C. K. Hall, J. Chem. Phys. **107**, 10691 (1997).

[33] K. Huang, *Statistical Mechanics* (Wiley, New York, 1963).