

## Reaction Path Optimization with Holonomic Constraints and Kinetic Energy Potentials

Jason B. Brokaw,<sup>†</sup> Kevin R. Haas,<sup>‡</sup> and Jih-Wei Chu<sup>\*,‡</sup>

*Department of Chemistry and Department of Chemical Engineering,  
University of California, Berkeley, California 94720*

Received March 24, 2009

**Abstract:** Two methods are developed to enhance the stability, efficiency, and robustness of reaction path optimization using a chain of replicas. First, distances between replicas are kept equal during path optimization via holonomic constraints. Finding a reaction path is, thus, transformed into a constrained optimization problem. This approach avoids force projections for finding minimum energy paths (MEPs), and fast-converging schemes such as quasi-Newton methods can be readily applied. Second, we define a new objective function – the total Hamiltonian – for reaction path optimization, by combining the kinetic energy potential of each replica with its potential energy function. Minimizing the total Hamiltonian of a chain determines a minimum Hamiltonian path (MHP). If the distances between replicas are kept equal and a consistent force constant is used, then the kinetic energy potentials of all replicas have the same value. The MHP in this case is the most probable isokinetic path. Our results indicate that low-temperature kinetic energy potentials (<5 K) can be used to prevent the development of kinks during path optimization and can significantly reduce the required steps of minimization by 2–3 times without causing noticeable differences between a MHP and MEP. These methods are applied to three test cases, the C<sub>7eq</sub>-to-C<sub>ax</sub> isomerization of an alanine dipeptide, the <sup>4</sup>C<sub>1</sub>-to-<sup>1</sup>C<sub>4</sub> transition of an α-D-glucopyranose, and the helix-to-sheet transition of a GNNQQNY heptapeptide. By applying the methods developed in this work, convergence of reaction path optimization can be achieved for these complex transitions, involving full atomic details and a large number of replicas (>100). For the case of helix-to-sheet transition, we identify pathways whose energy barriers are consistent with experimental measurements. Further, we develop a method based on the work energy theorem to quantify the accuracy of reaction paths and to determine whether the atoms used to define a path are enough to provide quantitative estimation of energy barriers.

### Introduction

Characterizing rare events in molecular systems is of critical importance in many fields of chemistry, material science, and biology. A large amount of effort has, thus, been devoted to developing computational methods to assist in identifying transition states and reaction paths.<sup>1–21</sup> For transitions between two metastable states, a simple and powerful

strategy is to employ a chain of replicas of the molecular system to connect one state to another.<sup>2,3</sup> A reaction path is obtained by optimizing an objective function of the chain. For example, minimizing the total potential energy of a chain determines a minimum energy path (MEP).<sup>4</sup> Other commonly employed objective functions for path optimization will be discussed later. Many computational methods such as elastic band,<sup>2,3</sup> nudged elastic band (NEB),<sup>5–9</sup> max-flux path,<sup>10,11</sup> action-derived molecular dynamics,<sup>12</sup> string,<sup>13–19</sup> harmonic Fourier bead,<sup>20</sup> and replica path<sup>21</sup> are based on the chain-of-states framework. Significant progress has been

\* Corresponding author. E-mail: jwchu@berkeley.edu.

<sup>†</sup> Department of Chemistry.

<sup>‡</sup> Department of Chemical Engineering.

made in different fields by applying these methods to understand the mechanisms of important transition processes.<sup>22,23</sup>

The major challenge of chain-of-states methods is ensuring chain continuity while optimizing the specified objective function, and this nature of dual objectives severely limits the stability, efficiency, and robustness of reaction path optimization.<sup>5–21</sup> Problems typically encountered include the overestimation of transition barriers and conditional convergence of reaction path optimization, especially when studying complex transitions. These issues also make it difficult to quantify the accuracy of reaction path optimization. In this work, we develop systematic approaches to overcome the difficulties resulting from the dual objectives of chain-of-states methods to broaden the applications of the chain-of-states methods. Two methodologies are developed: (1) use holonomic constraints to ensure equal distances between replicas and (2) use kinetic energy potentials to prevent the development of kinks during path optimization. Furthermore, we develop a method to quantify the accuracy of a reaction path by employing the work energy theorem. Our results indicate that the developments of this work significantly enhance the stability, efficiency, and robustness of reaction path optimization using a chain of replicas and enable applications to complex molecular systems. In the following, difficulties due to the dual objectives of chain-of-states methods will first be summarized. We then describe the two strategies for overcoming these difficulties and our approach of quantifying the accuracy of reaction path optimization.

To ensure chain continuity, ways to control the spacing between replicas along a chain are required.<sup>2,3</sup> Without loss of generality, we focus on maintaining equal distances between replicas. A commonly employed strategy is to add potentials that restrain the distances between replicas to the same value.<sup>2,3</sup> However, optimizing the combined objective function results in solutions that satisfy neither the equal-spacing requirements nor the criteria of optimizing an objective function.<sup>5,6</sup> In cases of optimizing MEPs, gradients of potential energy functions prevent restraint potentials from equally spacing replicas and cause sliding-down; gradients of restraint potentials prevent replicas from reaching minima on the hyperplanes perpendicular to a path and cause corner-cutting.<sup>5,6</sup>

To decouple the effects of the objective function and restraint potentials, a nudged elastic band (NEB) method was invented to project away force components that cause sliding-down and corner-cutting based on an estimation of tangent vectors.<sup>5,6</sup> Although a path with stationary NEB forces is indeed a MEP, projected forces are not conservative.<sup>5,8</sup> As a result, it is difficult to apply faster converging quasi-Newton methods for path optimization.<sup>8,24</sup> A formal strategy is to separate the dual objectives of equal spacing and optimizing the objective function when applying quasi-Newton methods,<sup>8</sup> but this approach makes the underlying optimization problem underdetermined. An alternative strategy is to employ the magnitude of projected forces as the objective function,<sup>25,26</sup> but this requires the Hessian matrices of replicas for computing a gradient vector and is, thus, limited to relatively small molecular systems.

Similar to NEB, the zero-temperature string (ZTS) method employs projected forces for finding a MEP.<sup>15</sup> Recently, a simplified string method was proposed that does not require explicit force projections.<sup>17</sup> Unprojected forces, i.e., directions of steepest descent, are used to update the configurations of replicas; sliding-down, caused by force components parallel to the path, is then removed by reparametrization.<sup>17</sup> Reparametrization is based on a continuous representation of a curve from a finite number of points via numerical interpolation and is a major difference between ZTS<sup>15,17</sup> and NEB;<sup>5,6</sup> in NEB, restraint potentials are used to control distances between replicas. As such, ZTS and the simplified string method both exhibit linear convergence as a steepest descent optimization.<sup>17</sup>

In summary, the dual objectives of path optimization, i.e., finding an optimal solution of the specified objective function and ensuring chain continuity, make it difficult to apply fast-converging quasi-Newton schemes for path optimization. This difficulty has limited the stability, efficiency, and robustness of chain-of-state methods, especially when a large number of degrees of freedom and a rugged potential energy surface are involved. To overcome these challenges, a plausible strategy is to enforce the equal-spacing requirements as holonomic constraints during path optimization. In this way, analytical theory can be precisely implemented in the form of constrained optimization without imposing ad hoc procedures such as adding restraint potentials or reparametrization. As a result, force projections are no longer required, and the Lagrange multipliers are used to balance the constraint forces and true forces with high numerical precision. With equal distance holonomic constraints, quasi-Newton methods and other optimization schemes that assume conservative forces can, thus, be readily applied. Parameterization of a continuous curve from a finite number of replicas also becomes unnecessary, thus, eliminating the dependence of path optimization on interpolation schemes.<sup>13,14</sup> Using holonomic constraints to control the distances between replicas also allows a straightforward implementation of general definitions of distances such as via a noncommutative rms (root-mean-squared) best-fit procedure,<sup>27</sup> which is particularly important for simulating transition processes of macromolecules.<sup>8,21</sup> In this work, we devise a simple, fast, and stable scheme to equally space replicas along a chain as constraints.

In addition to maintaining equal distances, ways to prevent the development of kinks (high curvatures) are also important for stable path optimization.<sup>13–19</sup> The tendency of developing kinks during path optimization increases with the number of replicas in a chain as well as the magnitude of forces involved in a transition.<sup>6</sup> For methods that rely on force projections, this intrinsic instability also depends on the ways tangent vectors are estimated and can be improved by adopting a definition based on the relative energies between replicas.<sup>6</sup> Another cause of highly curved paths is the underlying ruggedness of the potential energy (or free energy) surface. Ways of controlling curvatures along a path are important for stable and efficient path optimization. Commonly employed approaches include using angles or self-avoiding potentials or applying smoothing procedures.

Preventing kinks during path optimization introduces additional ad hoc components that may bias the results; therefore, physics based approaches are highly desired. We propose a new strategy of using potentials that represent the kinetic energy of each replica to prevent the development of kinks during path optimization. Kinetic energy potentials have a functional form that is quadratic in the mass-weighted distance between two replicas with a force constant, corresponding to the inverse of squared time. Combining kinetic energy potentials with the potential energy functions of replicas defines the total Hamiltonian of a chain, and minimizing this objective function determines a minimum Hamiltonian path (MHP). To our knowledge, this is the first time that a Hamiltonian objective function and MHP have been introduced. Different ways of setting the force constants of kinetic energy potentials can be devised. In this work, we consider a simple strategy of using a constant value for all kinetic energy potentials. By keeping the same distance between replicas, this approach ensures a constant kinetic energy along a chain and gives an isokinetic path. A fixed value of kinetic energy can also be maintained during path optimization by rescaling the force constants. A convenient way for estimating the magnitude of a kinetic energy potential is through the value of the corresponding temperature. With zero kinetic energy potentials, the objective function becomes the total potential energy, and a MEP is, thus, a 0 K MHP. In the high-temperature limit, kinetic energy potentials ensure that a straight line connecting two structures is the result of path optimization. Therefore, employing kinetic energy potentials provides a physically based approach to prevent the development of kinks during path optimization. For optimizing MEPs using low-temperature kinetic energy potentials, our results indicate that the convergence of path optimization is accelerated by 2–3 times.

In the rest of this paper, the details of equal distance holonomic constraints and kinetic energy potentials are presented. To illustrate the enhanced stability, efficiency, and robustness of path optimization, methods developed in this work are applied to find the reaction paths of three molecular transitions, the C<sub>7eq</sub>-to-C<sub>ax</sub> isomerization of an alanine dipeptide, the <sup>4</sup>C<sub>1</sub>-to-<sup>1</sup>C<sub>4</sub> transition of an α-D-glucopyranose, and the helix-to-sheet transition of a GNNQQNY heptapeptide. Finally, the conclusion is stated.

## Methodology

**A. Distance between Replicas.** For a molecular system composed of  $N$  atoms, a configuration state is specified by the three-dimensional position vectors,  $r_i$  ( $i = 1, N$ ), of each atom. A general way for defining the distance ( $\Delta l^{IJ}$ ) between two states (replicas)  $I$  and  $J$  is<sup>8,21</sup>

$$\Delta l^{IJ} = \sqrt{\frac{\sum_{i=1}^N w_i (r_i^I - \mathbf{U}^{IJ} r_i^J)^2}{\sum_{i=1}^N w_i}} \quad (1)$$

In eq 1,  $w_i$  is the weighting factor associated with each atom for measuring distance. If  $w_i$  corresponds to the mass of atom

$i$ , eq 1 determines the mass-weighted distance between configurations  $I$  and  $J$ . For atoms that are not directly involved in the process of interest, such as solvent molecules far away from the active site of an enzyme,  $w_i$  may be set to zero.<sup>8,21</sup> When computing distances between molecular structures, a common practice is to exclude the contribution from rigid body rotation and translation. The contribution from rigid body translation can be removed by putting the weighted center at the origin; to remove the contribution from rigid body rotation, we adopt an rms best-fit procedure to calculate a three-by-three unitary rotation matrix ( $\mathbf{U}^{IJ}$  in eq 1) that minimizes the distance between  $I$  and  $J$  using a Lagrangian.<sup>27</sup> The rms best-fit procedure is widely used for measuring structural differences between biomolecules; it has also been employed in the replica path method for reaction path optimization.<sup>8,21</sup>

**B. The Chain-of-States Framework and Equal Distance Constraints.** Given two metastable states of a molecular system of  $N$  atoms with their configurations denoted as  $\mathbf{r}^0 = \{r_i^0, i = 1, N\}$  and  $\mathbf{r}^K = \{r_i^K, i = 1, N\}$ , a chain of  $K + 1$  replicas can be constructed to connect the two states. To ensure chain continuity and resolution of path, we propose to constrain the distances between replicas being equal during path optimization, i.e.,

$$\Delta l^0 = \cdots = \Delta l^I = \cdots = \Delta l^{K-1} = \overline{\Delta l} \quad (2)$$

Here,  $\Delta l^I$  is a shorthand notation for  $\Delta l^{I,I+1}$  in eq 1, and  $\overline{\Delta l}$  is the averaged distance between replicas along a chain and is computed on the fly during path optimization. With  $K + 1$  replicas in a chain, eq 2 involves a set of  $K$  coupled algebraic equations that need to be solved at each step of path optimization. Following the procedures of maintaining holonomic constraints in molecular simulations,<sup>28</sup> we propose the following scheme for solving eq 2 with the two ends,  $\mathbf{r}^0$  and  $\mathbf{r}^K$ , fixed at the current position:

(i) Calculate the averaged distance,  $\overline{\Delta l}$ , between replicas from  $\{(\mathbf{r}^0)^{(0)} \cdots (\mathbf{r}^K)^{(0)}\}$ . The superscript “(n)” specifies the index of equal distance iteration.

(ii) Use a set of  $K$  coefficients,  $(\lambda^I)^{(n)} (I = 0, K - 1)$ , to update the position vector of each replica  $I$ :

$$(\mathbf{r}^I)^{(n+1)} = (\mathbf{r}^I)^{(n)} + (\lambda^{I-1})^{(n)} \left( \frac{\partial \Delta l^{I-1}}{\partial \mathbf{r}^I} \right)^{(n)} + (\lambda^I)^{(n)} \left( \frac{\partial \Delta l^I}{\partial \mathbf{r}^I} \right)^{(n)} \quad (3)$$

These coefficients are going to be solved in step (iii). For a chain of  $K + 1$  replicas, there are  $K$  segments along a chain, and  $K$  coefficients ( $I = 0, K - 1$ ) are involved. Since we fix the two ends, eq 3 is applied from replicas  $I$  to replica  $K - I$ .

(iii) Solve  $(\lambda^I)^{(n)}$  by setting the first order Taylor expansion of each of  $((\Delta l^I)^{(n+1)} - \overline{\Delta l})$  ( $I = 0, K - 1$ ) with respect to  $(\lambda^J)^{(n)}$  to zero, i.e.,

$$-((\Delta l^I)^{(n)} - \overline{\Delta l}) = \sum_{J=I-1}^{I+1} \left( \frac{\partial \Delta l^I}{\partial \lambda^J} \right)_{\lambda^J=0} (\lambda^J)^{(n)} \quad (4)$$

There are a total of  $K$  equations for  $K$  unknowns. Since each replica is only coupled to its nearest neighbors through

distance constraints, eq 4 is a set of tridiagonal algebraic equations, which can be solved using well-established methods.<sup>29</sup>

(iv) If any of the values of  $l((\Delta l^I)^{(n+1)} - \bar{\Delta l})$  ( $I = 0, K - 1$ ) calculated from  $\{(\mathbf{r}^0)^{(n+1)} \cdots (\mathbf{r}^K)^{(n+1)}\}$  (via eq. 3) is greater than tolerance, then repeat (ii–iii). A tolerance of  $10^{-8}$  Å is used throughout this work.

During the iterations for solving eq 2,  $\bar{\Delta l}$  remains fixed and the total length of a chain is left the same by steps i–iv;  $\bar{\Delta l}$ , though, can change freely during path optimization. Gradients of distances are employed in eq 3 as the basis set for adjusting the position of each replica, and  $\lambda^I$  designates the displacements along these vectors. Equal distance constraints are usually satisfied via steps ii–iv within 20–25 iterations, even during the initial stages of path optimization where large steps are often involved. The proposed scheme works equally well if the distances between replicas are defined via rms best-fit and demonstrate similar stability and efficiency for all molecular systems that we have tested so far. Generality to rms best-fit distance is a major advantage of using holonomic constraints instead of reparametrization schemes used in the string method. Compared to the use of restraint potentials, as in NEB, use of holonomic constraints guarantees high accuracy in the satisfaction of constraint equations via Lagrange's multipliers, separating it from the optimization of the objective function. Furthermore, this approach can be generalized if variable density of replicas along a path is required by a given mechanism of distributing distances, for example, using more replicas around higher energy regions.

**C. Objective Functions for Path Optimization.** The proposed procedure for solving distance constraints can be applied with any objective function for path optimization. Several commonly employed functions for path optimization are briefly summarized, followed by the introduction of kinetic energy potentials and minimum Hamiltonian paths.

**1. Minimum Energy Path (MEP).** A solution of minimizing the total potential energy,  $\sum_{I=0}^K V^I$ , of a chain subject to equal distance constraints (eq 2) determines a MEP<sup>4</sup> that satisfies:

$$(-\nabla_I V^I)(1 - \boldsymbol{\tau}^I \boldsymbol{\tau}^I) = 0 \quad (5)$$

$V^I = V(\mathbf{r}^I)$  is the potential energy of replica  $I$ , and  $\boldsymbol{\tau}^I$  is the unit tangent vector of replica  $I$ . Under the framework of constrained optimization,<sup>29</sup>  $\boldsymbol{\tau}^I$  are defined by the gradients of the constraint equations (eq 2) that involve replica  $I$ :

$$\boldsymbol{\tau}^I = \frac{\frac{\partial \Delta l^{I-1}}{\partial \mathbf{r}^I} + \frac{\partial \Delta l^I}{\partial \mathbf{r}^I}}{\left| \frac{\partial \Delta l^{I-1}}{\partial \mathbf{r}^I} + \frac{\partial \Delta l^I}{\partial \mathbf{r}^I} \right|} = \frac{\sum_{i=1}^N w_i (\mathbf{r}_i^{I+1} - \mathbf{r}_i^{I-1})}{\left| \sum_{i=1}^N w_i (\mathbf{r}_i^{I+1} - \mathbf{r}_i^{I-1}) \right|} \quad (6)$$

**2. Minimum Free Energy Path (MFEP).** Following Maragliano et al.,<sup>18</sup> a MFEP is defined as a path on which the mean force at a finite temperature is parallel to the path, i.e.,

$$\langle -\nabla_I F^I \rangle_I (1 - \boldsymbol{\tau}^I \boldsymbol{\tau}^I) = 0 \quad (7)$$

$F^I$  is the free energy for constraining the system at the configuration of replica  $I$ . Methods for optimizing MEPs can, thus, be generalized at a finite temperature to identify MFEPs,<sup>18</sup> and the approach of using holonomic constraints can be readily applied. Details of identifying the MFEP with equal distance constraints and computing free energy profiles along a path will be discussed in future work.

**3. Dynamic Path.** Dynamic paths can be found by applying the least action principle<sup>30</sup> where the kinetic energy associated with each segment of a path can be defined as:

$$T^I = \sum_{i=1}^N \frac{1}{2} m_i \left( \frac{r_i^{I+1} - r_i^I}{\Delta t^I} \right)^2 = \frac{1}{2} \frac{M}{(\Delta t^I)^2} (\Delta l_m^I)^2 = \frac{1}{2} M (v_m^I)^2 \quad (8)$$

In eq 8,  $M$  is the total mass of the molecular system and  $\Delta l_m^I$  is the mass-weighted distance between replica  $I$  and  $I + 1$  (see eq 1).  $\Delta t^I$  represents the time for the system to travel a distance of  $\Delta l_m^I$ , and  $\Delta l_m^I / \Delta t^I$  is the mass-averaged root-mean-squared velocity,  $v_m^I$ . The direction of velocity vector is specified by the difference between two position vectors. Given a value of  $\Delta t^I$ ,  $T^I$  gives the kinetic energy for connecting two neighboring replicas separated by  $\Delta l_m^I$ . Equation 8 is of the form of a potential energy function that is quadratic in  $\Delta l_m^I$  with  $M/(\Delta t^I)^2$  as the force constant; therefore,  $T^I$  is referred to as kinetic energy potential. Assigning a force constant to a kinetic energy potential is equivalent to assigning a value of the time duration for connecting two replicas. If all of the atoms in a molecular system are included in defining the distance between replicas (eq 1), then  $T^I$  corresponds to the total kinetic energy of the system excluding translational and rotational degrees of freedom. For cases where a subset of the system is used to define the path,  $T^I$  corresponds to the kinetic energy of the subset.

When the distances between replicas are constrained to be equal, it is natural to use the total contour length,  $L$ ,  $L = \sum_{I=0}^{K-1} \Delta l_m^I$ , instead of time as the independent variable for classical action.<sup>25,30</sup> A discrete representation of action  $S_L$  for a chain with  $K + 1$  replicas is<sup>25,30</sup>

$$S_L = \int_0^L d\ell \sqrt{2(E - V)} \approx \Delta l_m \sum_{I=0}^K v_m^I = \frac{L}{K} \sum_{I=0}^K v_m^I \quad (9)$$

A chain of replicas that corresponds to a stationary point of  $S_L$  satisfies  $(\delta S_L)/(\delta \mathbf{r}^I) = 0$ , which leads to Newton's equations of motion:<sup>30</sup>

$$\left[ \frac{m_i}{(\Delta t^I)^2} (r_i^{I+1} - 2r_i^I + r_i^{I-1}) \right] = [-\nabla_I V^I (1 - \boldsymbol{\tau}_m^I \boldsymbol{\tau}_m^I)]_i \quad (10)$$

In eq 10,  $\boldsymbol{\tau}_m^I$  is the tangent vector associated with replica  $I$ ; the subscript  $m$  indicates that atomic masses are the weighting factors for defining distances between replicas (see eq 6). We also implicitly assume that all of the atoms in the molecular system are involved in defining the distances between replicas. The values of  $v_m^I$  in eq 9 may be determined by the work–energy theorem, and  $\Delta t^I$  in eq 10 can then be calculated as  $\Delta l_m / v_m^I$ . Finding stationary solutions of eq 9,



however, is very difficult,<sup>12,25</sup> since they may include the minimum, maximum, or saddle points of  $S_L$ . Directly minimizing  $|\delta S_L|/(\delta \mathbf{r}^I)$  requires the calculations of the Hessian matrices of replicas and is limited to relatively small molecular systems.<sup>12</sup>

4. *Minimum Hamiltonian Path (MHP)*. The sum of the kinetic energy potential (eq 8) and potential energy of a replica defines its Hamiltonian,<sup>30</sup> and the total Hamiltonian of a chain is

$$H^{\text{tot}} = \sum_{I=0}^{K-1} (T^I + V^I) = \sum_{I=0}^{K-1} H^I \quad (11)$$

Kinetic-energy potentials characterize the magnitude of inertia for keeping a path straight, and thus, providing a physically based approach for reducing curvatures along a path. If path optimization is conducted only with  $T^I$ , then a straight line connecting two configurations would result. The kinetic energy potential  $T^I$  in eq 8 is quadratic in  $\Delta l_m^I$  with  $M/(\Delta t^I)^2$  as the force constant. Minimizing the total Hamiltonian of a chain subject to equal distance constraints gives a minimum Hamiltonian path (MHP) that satisfies:

$$\left[ -\frac{m_i}{(\Delta t^I)^2} (r_i^{I+1} - 2r_i^I + r_i^{I-1}) \right] = [-\nabla_I V^I (1 - \boldsymbol{\tau}_m^I \boldsymbol{\tau}_m^I)]_i \quad (12)$$

The left-hand side of eq 12 is the same as that of eq 10, but with an opposite sign. Therefore, a MHP does not satisfy Newton's equation of motion, but such dynamics could be achieved if the system is coupled to an external bath. With all  $\Delta l_m^I$ 's kept at the same value, assigning a common value to all force constants gives constant kinetic energy potentials along a path, i.e., an isokinetic path.<sup>31</sup> In this case, a MHP is also a minimum energy isokinetic path. Analogous to a MEP being the most probable path for an overdamped system,<sup>32</sup> a minimum energy isokinetic path is most highly likely among paths that satisfy isokinetic conditions.<sup>31</sup>

The temperature associated with a kinetic energy potential can be calculated to characterize its magnitude.<sup>31</sup> By scaling the force constant of kinetic energy potentials during path optimization, one may fix kinetic energy potentials at a desired value (or temperature) during path optimization. If kinetic energy potentials are set to zero, then the objective function becomes the total potential energy; a minimum energy path (MEP) is thus a 0 K MHP. Kinetic energy potentials also provide an estimation of the time scale of progression,  $t^I = \Delta l_m / v_m^I$ , associated with each segment. A low kinetic energy potential corresponds to a longer progression time. For a MEP,  $v_m^I \rightarrow 0$  and  $t^I \rightarrow \infty$ .

Using kinetic energy potentials is, thus, a physically based approach to control curvature during path optimization. Highly kinked paths can result from having too many replicas<sup>6</sup> or are due to the nature of objective functions. For example, when optimizing MFEP using the finite temperature string method, finite time molecular dynamics or Langevin dynamics are performed to calculate mean forces.<sup>16</sup> As a result, statistical noise is inevitably involved and can cause instability and kinks during path optimization.<sup>16</sup> In both cases, kinetic energy potentials can be used to maintain a

smooth path and stability during optimization. As shown later, the magnitude of kinetic energy potentials can be quantitatively controlled to minimize the effects on objective function.

**D. Quantifying Accuracy of Reaction Path Optimization.** On a MEP or MHP with isokinetics, the potential energy difference between replicas can be computed by applying the work–energy theorem, for which a discrete representation based on a second-order symmetric finite difference is

$$(\Delta V^I)^{\text{WET}} = \frac{1}{4} \sum_{J=I}^{I+1} (\nabla_J V^J)(\boldsymbol{\tau}_m^J)(\Delta l_m^{J+1,J-1}) \quad (13)$$

In eq 13, the length, over which  $(-\nabla_I V^I)(\boldsymbol{\tau}_m^I)$  is exerted, is approximated by  $\Delta l_m^{I+1,I-1}/2$ . Equation 13 or, equivalently, the work–energy theorem, can be used to quantify the accuracy of reaction path optimization. For example, accuracy can be measured by the rms differences between  $(\Delta V^I)^{\text{WET}}$  calculated via eq 13 and the calculated values of energy difference,  $\Delta V^I = (V^{I+1} - V^I)$ ,  $\delta(\Delta V)_{\text{rms}}$ , as

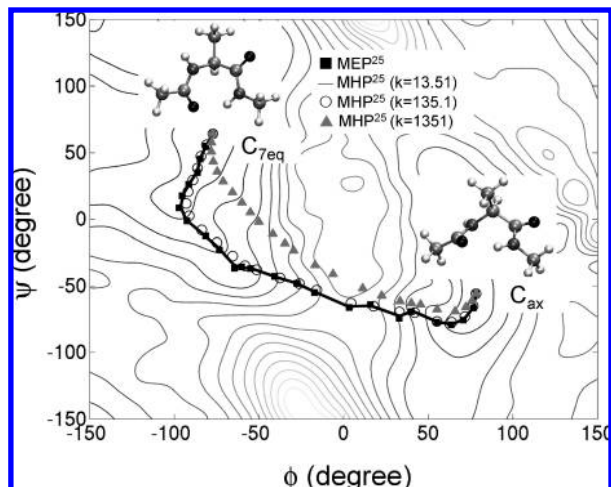
$$\delta(\Delta V)_{\text{rms}} = \sqrt{\frac{\sum_{I=0}^{K-1} ((\Delta V^I)^{\text{WET}} - \Delta V^I)^2}{K}} \quad (14)$$

If a reaction path is accurate, then the resulting path will give a small value of  $\delta(\Delta V)_{\text{rms}}$ . If the algorithm of reaction path optimization is correct, then the value of  $\delta(\Delta V)_{\text{rms}}$  is expected to decrease with the increase of the number of replicas used in a chain. As suggested in eq 13, factors that affect the value of  $\delta(\Delta V)_{\text{rms}}$  include the degrees of freedom used in defining distances between replicas, the estimation of tangent vectors, and the calculation of distances between replicas. Since reaction path methods may differ in these properties, eqs 13 and 14 provide an unbiased and quantitative measure of accuracy.

## Results and Discussion

To illustrate the enhancement in stability, efficiency, and robustness by employing equal distance holonomic constraints and kinetic energy potentials, they are applied to the path optimization of three molecular transitions: the C<sub>7eq</sub>-to-C<sub>ax</sub> transition of an alanine dipeptide, the <sup>4</sup>C<sub>1</sub>-to-<sup>1</sup>C<sub>4</sub> transition of an  $\alpha$ -D-glucopyranose, and the helix-to-sheet transition of a GNNQQNY heptapeptide. Since all-atom force fields are used, at least 20 replicas are required to represent transitions in all three cases for accurate reaction optimization (as measured by eqs 13 and 14). As such, both methods are found to be essential for the convergence of reaction path optimization. All calculations are performed with the CHARMM program,<sup>33</sup> based on which (version c35a2) we implemented the new developments. Simulation details will be described in each test case.

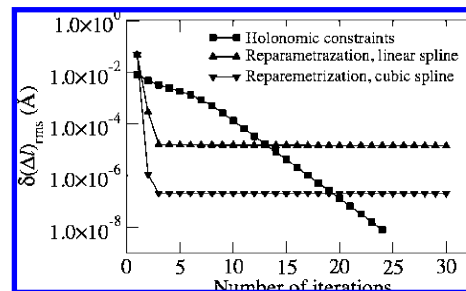
**A. Isomerization of Alanine Dipeptide.** The transition of alanine dipeptide (CH<sub>3</sub>–CO–NH–CHCH<sub>3</sub>–CO–NH–CH<sub>3</sub>) between C<sub>7eq</sub> ( $\phi = 77.3, \psi = 64.3$ ) and C<sub>ax</sub> ( $\phi = 78.2, \psi = -55.8$ ) isoforms (see Figure 1) is a standard test case for the development of computational methods that aim



**Figure 1.** Pathways of the  $C_{7eq}$ -to- $C_{ax}$  transition of an alanine dipeptide. Pathways are projected on a two-dimensional surface of backbone dihedral angles ( $\phi$ :  $C-N-C_{\alpha}-C$ ,  $\psi$ :  $N-C_{\alpha}-C-N$ ). The contour plot is generated via minimized structures of alanine dipeptide restrained at different values of  $\phi$  and  $\psi$ . Structures of the  $C_{7eq}$  and  $C_{ax}$  configurations are shown via a ball-and-stick representation. MEP stands for a minimum energy path, and MHP stands for a minimum Hamiltonian path; superscript indicates the number of replicas used in a chain. Values of the force constants of kinetic energy potentials ( $k$ 's) in kcal/mol/Å<sup>2</sup> are also specified. See text for discussion.

to find reaction paths. The CHARMM22 all-atom force field<sup>34</sup> with CMAP backbone dihedral angle corrections<sup>35</sup> is used for calculation without solvation, and all pair interactions between atoms are involved in calculating the potential energy. The initial configurations of replicas are generated by linearly interpolating in the  $\phi$  and  $\psi$  space.

All atoms except methyl hydrogen are used to define the distances between replicas using eq 1 with a rms best-fit procedure to remove contributions from rigid body translation and rotation using atomic masses as weighting factors. Although the initial configurations of replicas generated via linear interpolation are far from being equally spaced based on mass-weighted rmsd, the proposed procedure for solving the constraint equations still manages to converge within 25 iterations for a chain of 25 replicas. The criterion of convergence is that the root-mean-squared difference in rms distance ( $\delta(\Delta l)$ ) between replicas is less than  $10^{-8}$  Å. Using steps i–iv as described earlier in section B, we also found that the required number of iterations for converging equal distance constraints is not sensitive to the number of replicas used in the chain. Compared to using restraint potentials, which corresponds to a single step of iteration, employing holonomic constraints requires additional cost in solving constraint equations to gain accuracy and to decouple maintaining equal distances from optimizing the objective function. Another widely used approach to maintain equal distances between replicas is reparametrization using interpolation functions such as cubic splines.<sup>15,17</sup> If the rms best-fit procedure is not used to define the distance between replicas, then forming a continuous curve to connect replicas is straightforward and reparametrization is very efficient in equating the distances between replicas. Using the same

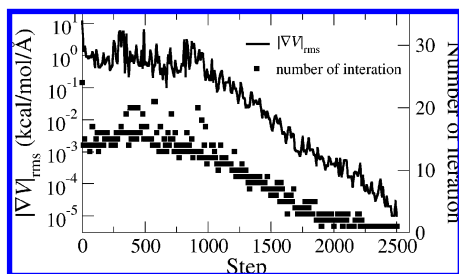


**Figure 2.** The rms difference in distances between replicas as a function of iteration for solving constraint equations and reparametrization. The initial path of the  $C_{7eq}$ -to- $C_{ax}$  transition using 25 replicas is used as the testing case. See text for details.

initial path as described above, Figure shows that only four iterations of reparametrization are needed to reduce  $\delta(\Delta l)$  to an accuracy determined by the interpolation functions.  $\delta(\Delta l)$  using cubic splines increases the upper limit of accuracy by two orders of magnitude ( $\delta(\Delta l) = 2 \times 10^{-6}$  Å) compared to that of using a linear interpolation ( $\delta(\Delta l) = 2 \times 10^{-4}$  Å). On the other hand, since holonomic constraints only require the input of constraint equations, the accuracy of convergence does not depend on the choice of interpolation functions. Furthermore, if a rms best-fit procedure is applied to calculating distances between replicas, then using the holonomic constraints approach allows straightforward generalization, whereas the applicability of continuous curve approximation in the reparametrization approach requires further investigation, which is beyond the scope of this work. In both cases of using reparametrization shown in Figure 2, the rms best-fit procedure is not used in calculating distances between replicas. In conclusion, using holonomic constraints provides higher accuracy in satisfying the equal distance constraints and higher flexibility in defining the distance between replicas, as compared to other commonly used methods. Additional cost, though, is involved in solving the constraint equations. As shown in the following, decoupling the optimization of the objective function and maintaining equal distances between replicas using holonomic constraints play important roles in enhancing the stability and ensuring the accuracy of reaction path optimization. Furthermore, the percentage cost of solving constraint equations would decrease with the increase of system size and the complexity of energy calculations, for example, when ab initio methods are used.

The rms gradient of potential energy,  $|\nabla V|_{rms}$ , for a chain of 25 replicas during the course of ABNR (adopted-basis Newton–Raphson) minimization is shown in Figure 3 (left axis, log scale); also shown (right axis) is the number of iterations for solving constraint equations. Kinetic energy potentials are not used here. The superlinear convergence of ABNR minimization is clearly seen in Figure 3, with a convergence criterion of  $|\nabla V|_{rms} < 10^{-5}$  kcal/mol/Å satisfied within 2 500 steps. By enforcing equal distance constraints, we illustrate that a MEP can be obtained without using restraint potentials, reparametrization, and/or force projections.

During the path optimization shown in Figure 3, the mass-weighted rmsd's ( $rmsd_M$ 's) between replicas are constrained

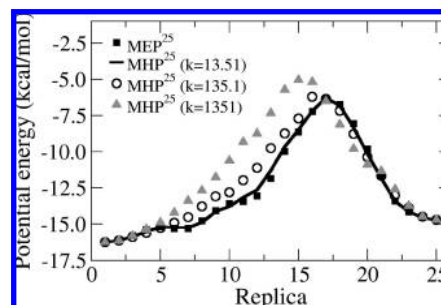


**Figure 3.** The rms gradient of potential energy  $|\nabla V|_{\text{rms}}$  (labels on left, y-axis) and the number of iterations for solving the constraint equations of equal distances (labels on right, y-axis) during path optimization using ABNR. The chain contains 25 replicas of alanine dipeptide to connect  $C_{7\text{eq}}$  and  $C_{\text{ax}}$ . Kinetic energy potentials are not used in path optimization.

to be equal. All atoms other than methyl hydrogen (13 out of 22) are used to define the distance between two replicas. As shown in Figure 3, the constraint equations can be satisfied within  $\sim 20$  iterations via steps i–iv. As the optimization approaches convergence, the number of iterations also decreases since the magnitude of displacement per step also decreases. The value of  $\text{rmsd}_M$  of the converged MEP with 25 replicas is  $0.103 \text{ \AA}$ , and the total contour length of the MEP is, thus,  $2.48 \text{ \AA}$ ; the  $\text{rmsd}_M$  between  $C_{7\text{eq}}$  and  $C_{\text{ax}}$  is  $1.17 \text{ \AA}$ . The 25 replica MEP ( $\text{MEP}^{25}$ ) is projected onto a  $(\phi, \psi)$  surface and shown in Figure 1. Without using holonomic constraints, reaction path optimization often fails to converge using restraint potentials alone. This tendency becomes more significant in the other two test cases since the underlying transition is more complicated.

Combining potential energy functions and kinetic energy potentials as the objective function of path optimization determines a minimum Hamiltonian path (MHP). MHPs of 25 replicas ( $\text{MHP}^{25}$ ) with different magnitudes of kinetic energy potentials are computed to illustrate their effects on the resulting path; results are shown in Figure 1. With a  $k \equiv M/(\Delta t)^2$  value of  $13.51 \text{ kcal/mol/\AA}^2$ , the resulting MHP (solid line in Figure 1) is indistinguishable compared to the MEP (solid squares). The  $\text{rmsd}_M$  between replicas on  $\text{MHP}^{25}$  ( $k = 13.51 \text{ kcal/mol/\AA}^2$ ) is  $0.1 \text{ \AA}$  (the value for the MEP is  $0.103 \text{ \AA}$ ), and the corresponding temperature of kinetic energy potentials has a very small value of  $1.74 \text{ K}$ . When  $k$  is increased to  $135.1 \text{ kcal/mol/\AA}^2$ , the resulting  $\text{rmsd}_M$  between replicas on the MHP becomes  $0.092 \text{ \AA}$  and the corresponding temperature is  $14.8 \text{ K}$ . The projected  $\text{MHP}^{25}$  ( $k = 135.1 \text{ kcal/mol/\AA}^2$ ) on a  $(\phi, \psi)$  surface (open circles in Figure 1) is also close to that of the MEP but is smoother near minima. When  $k$  is raised to  $1351 \text{ kcal/mol/\AA}^2$ , the  $\text{rmsd}_M$  between replicas becomes  $0.068 \text{ \AA}$  and the corresponding temperature is  $80.61 \text{ K}$ . From Figure 1, the deviation of  $\text{MHP}^{25}$  ( $k = 1351 \text{ kcal/mol/\AA}^2$ ; solid triangles) from the MEP is clear: the path is straightened by kinetic energy potentials near local minima but to a much lesser extent near a saddle point.

The pathways shown in Figure 1 indicate that kinetic energy potentials affect more on regions with a flat potential energy surface, over which kinks tend to develop during path optimization.<sup>6</sup> This trend can also be seen in the plots of potential energy profiles shown in Figure 4. Increasing the



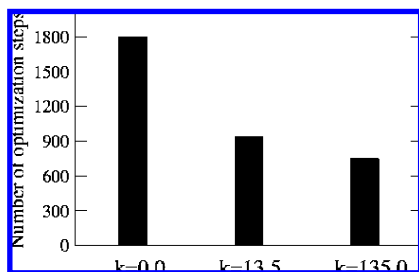
**Figure 4.** Profiles of potential energy on optimized pathways of the  $C_{7\text{eq}}$ -to- $C_{\text{ax}}$  transition of an alanine dipeptide in the gas phase. The x-axis specifies the index of the replica along a chain. MEP stands for a minimum energy path, and MHP stands for a minimum Hamiltonian path; superscript indicates the number of replicas used in a chain. Values of the force constants of kinetic energy potentials ( $k$ 's) in  $\text{kcal/mol/\AA}^2$  are also specified.

strength of kinetic energy potentials shifts the potential energy profile to the left, since the  $C_{7\text{eq}}$  basin (replica 1) is flatter than that of  $C_{\text{ax}}$ . The shoulder between replicas 5 and 8 on  $\text{MEP}^{25}$  also gradually disappears as  $k$  is increased from  $13.51$  to  $1351 \text{ kcal/mol/\AA}^2$ . The value of the highest potential energy along a path (replica 17 for  $\text{MEP}^{25}$ ,  $\text{MHP}^{25}$  ( $k = 13.51$  and  $135.1 \text{ kcal/mol/\AA}^2$ ) and replica 15 for  $\text{MHP}^{25}$  ( $k = 1351 \text{ kcal/mol/\AA}^2$ )) does increase due to the presence of kinetic energy potentials but only becomes appreciable when  $k = 1351 \text{ kcal/mol/\AA}^2$ .  $\Delta E^{\text{max}}$ , the difference in energy between the highest energy replica on  $\text{MHP}^{25}$  and  $\text{MEP}^{25}$  is  $1.31 \text{ kcal/mol}$  for  $k = 1351 \text{ kcal/mol/\AA}^2$  but is only  $0.07$  and  $0.001 \text{ kcal/mol}$  for  $k = 135.1$  and  $13.51 \text{ kcal/mol/\AA}^2$ , respectively. Differences in geometries due to kinetic energy potentials also follow a similar trend.  $\text{Rmsd}_M^{\text{max}}$ , the  $\text{rmsd}_M$  between the highest energy replica on  $\text{MHP}^{25}$  and  $\text{MEP}^{25}$ , is  $0.2$  for  $k = 1351$ ,  $0.04$  for  $k = 135.1$ , and  $0.009 \text{ \AA}$  for  $k = 13.51 \text{ kcal/mol/\AA}^2$ .

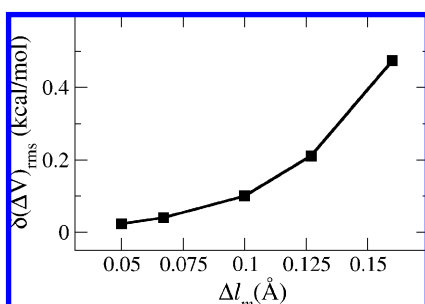
In summary, using 25 replicas to represent the  $C_{7\text{eq}}$ -to- $C_{\text{ax}}$  transition of an alanine dipeptide, kinetic energy potentials with  $k < 150 \text{ kcal/mol/\AA}^2$  do not cause an appreciable difference between the resulting MHP and MEP ( $\Delta E^{\text{max}} < 0.1 \text{ kcal/mol}$ ;  $\text{rmsd}_M^{\text{max}} < 0.1 \text{ \AA}$ ). On the other hand, the number of steps for converging reaction path optimization ( $|\nabla V|_{\text{rms}} < 10^{-5} \text{ kcal/mol/\AA}$ ) is much less (by 2–3 times) in the presence of kinetic energy potentials due to the reduced occurrence of kinks; results are shown in Figure 5. Even in the presence of a large number of replicas (up to 97), kinetic energy potentials that correspond to a temperature of  $\sim 2$ – $5 \text{ K}$  suffice to ensure stable path optimization. This level of kinetic energy potential ( $< 5 \text{ K}$ ) causes a negligible difference in energetics ( $\Delta E^{\text{max}} < 0.05 \text{ kcal/mol}$ ) and geometries ( $\text{rmsd}_M^{\text{max}} < 0.05 \text{ \AA}$ ) compared to a MEP. Without kinetic energy potentials, path optimization becomes difficult to converge when a large number of replicas are used in a chain.

Therefore, kinetic energy potentials can be used to significantly enhance the stability of reaction path optimization and increase the efficiency by reducing the steps of optimization. These advantages can be exploited without causing appreciable perturbation in structures and energetics. All calculations presented above are obtained by setting a





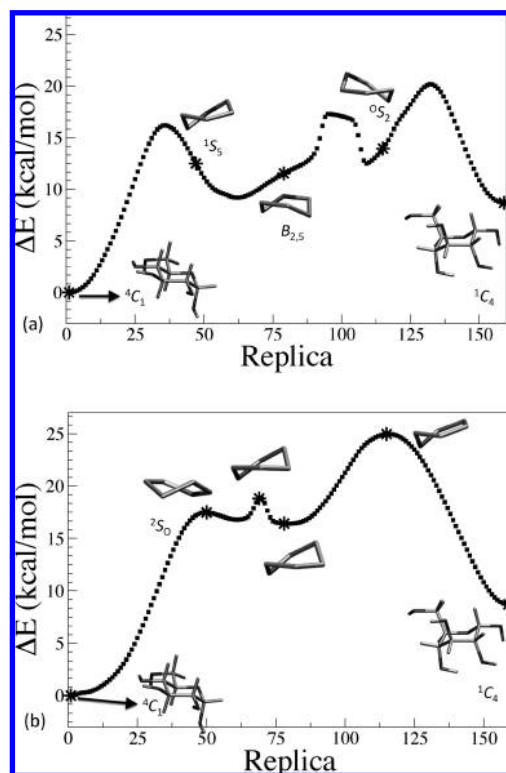
**Figure 5.** The required number of minimization steps for converging reaction path optimization ( $|\nabla V|_{\text{rms}} < 10^{-5}$  kcal/mol/Å) using different values of force constants for kinetic energy potentials. The  $C_{7\text{eq}}$ -to- $C_{\text{ax}}$  transition of an alanine dipeptide using 25 replicas is used as the example.



**Figure 6.** Rms error of using the work-energy theorem in computing potential energy differences between replicas via eq 13. Rms errors are plotted as a function of the  $\text{rmsd}_M$  between two neighboring replicas.

constant value of the force constants for kinetic energy potentials ( $k \equiv M/(\Delta t)^2$ , eq 8); an alternative approach is to fix the value of kinetic energy,  $1/2M(\Delta l/\Delta t)^2$ , or temperature by rescaling force constants during path optimization. When the same value of the force constant is used for all replicas and all distances between replicas are constrained to be equal, an isokinetic path results.<sup>31</sup> In this case, a MHP corresponds to a minimum energy isokinetic path, which is also the most probable path among pathways that satisfy isokinetic conditions.<sup>31</sup> Different schemes may also be designed to specify the force constants of kinetic energy potentials but will not be discussed in this work. Moreover, kinetic energy potentials can also be employed to systematically examine the effects of inertia on reaction pathways, as shown in Figure 1 and Figure 4.

A key parameter for chain-of-states methods is the number of states or replicas, and an important question is how many is considered enough. While it is difficult to predict a priori the required number of replicas for a transition, a useful metric is whether the work-energy theorem is satisfied by eqs 13 and 14. For reaction paths of the  $C_{7\text{eq}}$ -to- $C_{\text{ax}}$  transition of an alanine dipeptide with different numbers of replicas,  $\delta(\Delta V)_{\text{rms}}$  is plotted as a function of  $\Delta l_m$ , the mass-weighted rms best-fit distance between replicas, in Figure 6. It can be seen that for  $\Delta l_m = 0.1$  Å,  $\delta(\Delta V)_{\text{rms}} = 0.1$  kcal/mol. Since a second-order scheme is employed in eq 13 to estimate  $(\Delta V^j)_{\text{wet}}$ ,  $\delta(\Delta V)_{\text{rms}}$  decreases quadratically with  $\Delta l_m$  as shown in Figure 6. The small value of  $\delta(\Delta V)_{\text{rms}}$  for the case of alanine dipeptide also indicates that ignoring the hydrogen atoms of methyl groups in defining the distance between



**Figure 7.** Potential energy profiles of two optimized pathways, (a) and (b), of the  ${}^4C_1$ -to- ${}^1C_4$  transition of  $\alpha$ -D-glucopyranose in the gas phase. The x-axis specifies the index of the replica. MHPs are optimized by using kinetic energy potentials with a low temperature ( $<5$  K) for both paths. Structures of selected replicas are shown to illustrate the nature of the transition. See text for discussion.

replicas results in negligible error in computing the accumulated work along a reaction path.

To our knowledge, quantitative analysis of the accuracy of a reaction path via physically based approaches such as the work-energy theorem has not yet been conducted. Based on this criterion, the convergence of  $\delta(\Delta V)_{\text{rms}}$  with the reduction of  $\Delta l_m$  in Figure 6 justifies the validity of our scheme of reaction path optimization, and we found that both the use of rms best-fit procedure and holonomic constraints are critical in ensuring the accuracy of path optimization. As described earlier, another factor that affects  $\delta(\Delta V)_{\text{rms}}$  is the degrees of freedom used in defining distances between replicas and eqs 13 and 14 can be used to systematically examine whether an atom can be ignored in defining the reaction path or not. Figure 6 indicates that ignoring methyl hydrogen results in a negligible effect on the energy profile; ignoring any of the heavy atoms, though, results in a noticeable change ( $>1.0$  kcal/mol) in energy profile.

**B.  ${}^4C_1$ -to- ${}^1C_4$  Transition of  $\alpha$ -D-Glucopyranose.** D-Glucose is an important monosaccharide in medicine, energy storage, and biology, and  $\alpha$ -D-glucopyranose is an abundant anomer within this class of biomolecules. The properties and functions of polysaccharides are tightly related to the conformation of glucose monomers, and it is crucial to understand the relative energetics and interconversion between different conformational states.<sup>36–39</sup> As shown in Figure 7, the most stable conformation of  $\alpha$ -D-glucopyranose is the  ${}^4C_1$  chair, and it can undergo structural transition to a



less stable  ${}^1C_4$  chair. The  ${}^4C_1$ -to- ${}^1C_4$  transition involves collective rearrangements of atoms in the pyranose ring and side-chain hydroxyl and the methylhydroxyl groups and, thus, requires a robust reaction path method to analyze mechanisms. Holonomic equal distance constraints and kinetic energy potentials are applied to simulate the pathways of  ${}^4C_1$ -to- ${}^1C_4$  transition of an  $\alpha$ -D-glucopyranose molecule in the gas phase. Since a large magnitude of forces are involved in the transition, a large number of replicas is required. Both methods of employing holonomic constraints and kinetic energy potentials are found to be necessary for reaction path optimization to converge.

The CHARMM35 carbohydrate force field<sup>40</sup> is used to describe the potential energy surface of  $\alpha$ -D-glucopyranose conformation without using a cutoff for nonbonded interactions. The minimized *tg* isomers<sup>39,41</sup> of  ${}^4C_1$  and  ${}^1C_4$  are used as the initial and final replicas. The difference in energy between  ${}^4C_1$  and  ${}^1C_4$  is 8.7 kcal/mol, close to the result of 7.93 kcal/mol obtained from DFT calculations at the B3LYP/6-311++G\*\* level.<sup>42</sup> All atoms are used to define the distance between replicas via mass-weighted rmsd best-fit. Since a significant number of hydrogen atoms are involved, a large number of replicas are needed to resolve the structural arrangements during transition. From an initial path generated via linear interpolation, a path of 159 replicas is optimized and shown in Figure 7(a). The  $\text{rmsd}_M$  between replicas is 0.02 Å, and the magnitude of kinetic energy potentials used for path optimization is 3 K. A small  $\text{rmsd}_M$  of 0.02 Å between replicas ensures a small  $\delta(\Delta V)_{\text{rms}}$  of 0.04 kcal/mol.

Based on the conformational nomenclature of six-membered rings,<sup>41</sup> the path in Figure 7(a) follows a route of  ${}^4C_1$ - ${}^1S_5$ - $B_{2,5}$ - ${}^0S_2$ - ${}^1C_4$  (*S*: skew; *B*: boat), and the structures of selected conformations are highlighted in Figure 7(a). Although rotations of side-chain atoms are also involved, the first and third barriers in Figure 7(a) mainly correspond to rearrangements of the six-member ring. The highest energy state is 20.2 kcal/mol higher in energy compared to  ${}^4C_1$ , consistent with the results of DFT calculations.<sup>37</sup> The second barrier of Figure 7(a), on the other hand, is associated with the rearrangements of hydroxyl and methylhydroxyl groups without noticeable changes in the ring conformation. This result indicates that side-chain atoms are explicitly involved in the activation process of the  ${}^4C_1$ -to- ${}^1C_4$  transition.

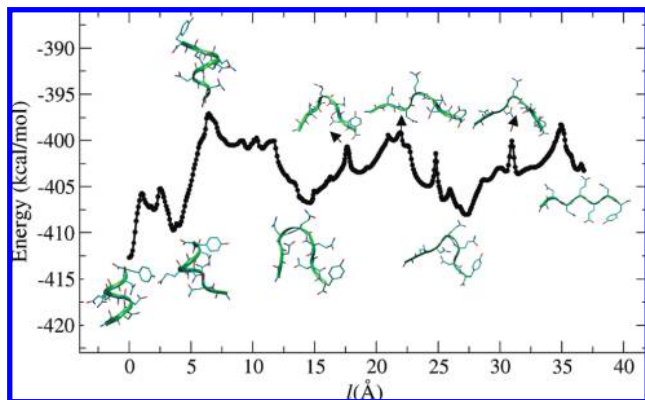
Since a  $B_{2,5}$ -like intermediate is found in the path shown in Figure 7a, we examine whether a path with  ${}^{2,5}B$  boat intermediate can also be identified. Using  ${}^{2,5}B$  boat as an intermediate for linear interpolation, an initial path is generated to test this mechanism. The optimized path with 159 replicas is shown in Figure 7b, with an  $\text{rmsd}_M$  of 0.02 Å between replicas. The magnitude of kinetic energy potentials used for path optimization is 3 K. Although replicas with a  ${}^2S_0$ -like structure are found along the path shown in Figure 7b (in contrast to the  ${}^0S_2$  structure in the path of Figure 7a), structures similar to  ${}^{2,5}B$  boat are not identified due to steric effects of the methylhydroxyl group. This result further affirms the importance of side-chain atoms in conformational transitions of glucopyranose. The path of Figure 7b has a higher barrier of 25.0 kcal/mol compared to that of the path in Figure 7a.

Combining holonomic equal distance constraints and kinetic energy potentials allows the identification of smooth pathways that use all 24 atoms of  $\alpha$ -D-glucopyranose to define distances between replicas. Three barrier-crossing events are involved in both pathways shown in Figure 7. Activation processes include rearrangements of side-chain atoms in addition to those of the pyranose ring. These results highlight the complexity of glucopyranose conformational changes and the requirement of a stable and efficient reaction path optimization scheme to analyze mechanisms. As such, this test case has not been widely attempted by other methods of reaction path optimization. Other pathways of the  ${}^4C_1$ -to- ${}^1C_4$  transition can be found by generating alternative initial paths via different interpolation strategies (currently ongoing work). Pathways obtained from molecular mechanical force fields may then be refined via *ab initio* methods.

**C. Helix-to-Sheet Transition of a GNNQQNY Heptapeptide.** The GNNQQNY heptapeptide is in an N-terminal prion-determining domain of Sup35, a prion-like protein in yeast that forms fibrillar amyloid assemblies.<sup>43</sup> Prions (proteinaceous infectious particles) are the cause of several neurodegenerative diseases and are formed by a prion protein whose conformation undergoes transition from the normal form to a misfolded form.<sup>44–46</sup> The GNNQQNY heptapeptide exhibits the amyloid properties of full-length Sup35,<sup>43</sup> and the helix-to-sheet transformation of this has crucial implications in the fibril formation of the full-length prion protein.<sup>44–46</sup>

The helix form of the GNNQQNY peptide is constructed according to the canonical  $\alpha$ -helical geometry, and the  $\beta$ -sheet form is obtained from X-ray crystallography (PDB code 1YJP).<sup>43</sup> The CHARMM22 all-atom protein force field with CMAP corrections was used with the FACTS implicit solvent model<sup>47</sup> to describe the potential energy surface of the GNNQQNY heptapeptide. Both  $\alpha$ -helix and  $\beta$ -sheet structures are stable local minima, with the  $\alpha$ -helix structure being 9.5 kcal/mol lower in potential energy. Including all atoms except the N-terminal hydrogen atoms to estimate the structural difference, the mass-averaged rms difference ( $\text{rmsd}_M$ ) between the  $\alpha$ -helix and  $\beta$ -sheet structures of the GNNQQNY peptide is 5.42 Å. It is expected that a pathway of the helix-to-sheet transition would involve multiple intermediate states and saddle points, thus, providing a stringent test for the stability and efficiency of reaction path optimization.

An initial path ( $\sim 20$  replicas) is first generated by restraining the molecule to different values of (( $\text{rmsd}_M$ -to-helix)-( $\text{rmsd}_M$ -to-sheet)) to connect the two structures. Even though the initial configurations are far from equally spaced, steps i–iv still manage to converge, indicating the robustness of the proposed scheme. As in the case of alanine dipeptide, an  $\text{rmsd}_M$  of 0.1 Å is needed to reduce the value of  $\delta(\Delta V)_{\text{rms}}$  (eq 14) to smaller than 0.1 kcal/mol and is used as the criterion for determining the number of replicas in a chain. Path optimization starts with higher kinetic energy potentials (path temperature  $\approx 15$  K), and their magnitude is then reduced to a temperature lower than 3 K, which ensures stable path optimization with minimal effects on the potential energy profile as compared to a MEP. If the  $\text{rmsd}_M$  of a converged path is larger than 0.1 Å, a new replica is inserted



**Figure 8.** Potential energy profiles of an optimized path of the helix-to-sheet transition of a solvated GNNQQNY heptapeptide as a function of contour length. The path contains a total of 464 replicas. Structures of selected replicas are shown to illustrate the nature of the transition. See text for discussion and the strategy for optimization.

in between each pair of replicas to double the number of segments until the  $\text{rmsd}_M$  between replicas of a converged path is around 0.1 Å.

Although steps i–iv can equally space a chain of hundreds of replicas (for the GNNQQNY peptide, the maximum number of replicas tested is 1 024), the number of steps required for path optimization also increases. Therefore, we use up to 100 replicas in a chain for path optimization. If doubling the number of segments along a chain results in more than 100 replicas, then the chain is divided into two segments at the location of a stable intermediate state. For the helix-to-sheet transition, this divide-and-conquer strategy results in a total of seven segments and a total of 456 replicas of the jointed chain. Potential energies of all replicas along the jointed chain and structures of selected intermediates and saddle points are shown in Figure 8.

For the path shown in Figure 8, the  $\alpha$ -helix starts to unwind from the N-terminal and then from the C-terminal end, along which the replica of the highest energy (15.5 kcal/mol higher than the  $\alpha$ -helix) is located. The peptide then goes through a series of configurations with coil and hairpin-like structures to reach to the  $\beta$ -sheet structure. Multiple intermediates and saddle points are identified; the total contour length (measured by  $\text{rmsd}_M$ ) for the chain is 36.7 Å,  $\sim 7$  times larger than the  $\text{rmsd}_M$  between the  $\alpha$ -helix and  $\beta$ -sheet structures (5.42 Å).

The computed barrier of 15.5 kcal/mol for the helix-to-sheet transition of the GNNQQNY peptide is lower than the activation energy (20–40 kcal/mol) of the helix-to-sheet transition of a longer segment (142-residue) of a mouse prion protein estimated by CD spectra at different temperatures.<sup>46</sup> By combining holonomic equal distance constraints and kinetic energy potentials with a divide-and-conquer strategy, we identify a pathway whose energetics are consistent with experimental measurements. This result has been difficult to achieve by applying chain-of-states methods<sup>46,48</sup> due to the limited stability and efficiency discussed earlier and demonstrates that the methods developed in this work can be used to enable the convergence of reaction path optimization of complex molecular systems.

For complicated transitions such as the helix-to-sheet transition, there may be multiple pathways to connect two states. While sampling reaction pathways is not the focus of this work, it is important to emphasize that holonomic equal distance constraints and kinetic energy potentials combined with the divide-and-conquer approach provide an efficient framework for exploring different initial paths and mechanisms. Initial paths can be generated using different strategies, for example, from fast pulling simulations along different order parameters or using different combinations of internal coordinates for interpolation. With the enhanced robustness and efficiency using the two new methods presented in this work, high-quality initial paths are not required for reaction path optimization. Alternative pathways can also be found by initiating molecular dynamics simulations from an optimized path, following the strategy of transition-path sampling.<sup>49–51</sup> Applications of these methods of helix-to-sheet transition of the GNNQQNY peptide and other conformational changes of biomolecules are ongoing.

## Conclusion

Two new methods are developed in this work to enable and accelerate the convergence of reaction path optimization using a chain of replicas, especially for transitions occurring in complex molecular systems. First, the requirements of maintaining equal distances between replicas are achieved by holonomic constraints to transform the finding of a reaction path to a problem of constrained optimization. This approach allows precise implementation of analytical theory without imposing ad hoc procedures such as adding restraint potentials or reparametrization. For finding a minimum energy or minimum free energy path, force projections are not required and fast-converging optimization schemes such as quasi-Newton methods can be readily applied. A simple, fast, and robust scheme is developed for solving the constraint equations that demonstrates high stability and efficiency even if complex molecular systems and many replicas are involved. This scheme also supports the use of a rms best-fit procedure to measure the distance between replicas.

Second, we propose to use kinetic energy potentials to prevent the development of kinks during path optimization and introduce a new objective function, the total Hamiltonian of a chain of replicas, for reaction path optimization. Our results indicate that by reducing the development of kinks, the number of steps for optimizing a path can be reduced significantly, by 2–3 times. By using low-temperature kinetic energy potentials, this enhancement in stability and efficiency can be acquired without causing appreciable differences in structures and energetics between minimum energy and minimum Hamiltonian paths.

To quantify the accuracy of reaction path optimization, a method is developed based on the work–energy theorem (eqs 13 and 14). The results of this analysis (Figure 6) justify the validity of our scheme of reaction path optimization. To our knowledge, a quantitative assessment of the accuracy of reaction path optimization has not yet been conducted. The analysis proposed in this work has solid theoretical foundation and can be used to test different strategies and

hypothesis used in reaction path optimization, such as the degrees of freedom, that should be used for defining distances between replicas.

To illustrate the enhanced stability, efficiency, and robustness of reaction path optimization by applying the methods developed in this work, they are applied to three test cases that involve rugged potential energy surfaces and high dimensionality: the C<sub>7eq</sub>-to-C<sub>ax</sub> transition of alanine dipeptide, the <sup>4</sup>C<sub>1</sub>-to-<sup>1</sup>C<sub>4</sub> transition of an α-D-glucopyranose, and the helix-to-sheet transition of a GNNQQNY heptapeptide. All-atom force fields are used in all three cases. By employing holonomic constraints to maintain equal distances between replicas and kinetic energy potentials to prevent kinks, convergence of reaction path optimization can be achieved even when a large number of replicas (>100) are used. Furthermore, convergence of path optimization is superlinear and is achieved by using an adopted-basis Newton–Raphson (ABNR) method. In this case, the application of fast-converging optimization schemes is straightforward since force projections are not involved due the use of holonomic constraints. The enhanced efficiency of reaction path optimization allows the identification of pathways of complex molecular systems whose barriers of transition lie in physiologically relevant range and can be compared with experimental measurements.

Finally, we observe that the optimized pathways in all three test cases demonstrate rearrangements of distinct parts of a molecular system, indicating the difficulty of using a few simple order parameters to describe the underlying transition. Using a chain of replicas to define a one-dimensional curve appears to be a general framework to overcome this difficulty. The methods developed in this work are easy to implement in other chain-of-states methods and can be used to enable the applications of this strategy to a broader range of molecular systems.

**Acknowledgment.** This work is supported by the National Institutes of Health through a training grant (no. T32GM008352), National Science Foundation through the graduate research fellowship program, and the University of California, Berkeley.

## References

- (1) Miller, W. H.; Handy, N. C.; Adams, J. E. *J. Chem. Phys.* **1980**, *72*, 99–112.
- (2) Elber, R.; Karplus, M. *Chem. Phys. Lett.* **1987**, *139*, 375–380.
- (3) Czerminski, R.; Elber, R. *Int. J. Quantum Chem.* **1990**, *38*, 167–186.
- (4) Olender, R.; Elber, R. *J. Mol. Struct.* **1997**, *398*, 63–71.
- (5) Jónsson, H.; Mills, G.; Jacobsen, K. W. In *Classical and Quantum Dynamics in Condensed Phase Simulations*; Berne, B. J., Ciccotti, G., Coker, D. F., Eds.; World Scientific: Singapore, 1998; p 185.
- (6) Henkelman, G.; Jónsson, H. *J. Chem. Phys.* **2000**, *113*, 9978–9985.
- (7) Henkelman, G.; Uberuaga, B.; Jónsson, H. *J. Chem. Phys.* **2000**, *113*, 9901–9904.
- (8) Chu, J.; Trout, B.; Brooks, B. *J. Chem. Phys.* **2003**, *119*, 12708–12717.
- (9) Trygubenko, S.; Wales, D. *J. Chem. Phys.* **2004**, *120*, 2082–2094.
- (10) Crehuet, R.; Field, M. J. *J. Chem. Phys.* **2003**, *118*, 9563–9571.
- (11) Huo, S. H.; Straub, J. E. *J. Chem. Phys.* **1997**, *107*, 5000–5006.
- (12) Passerone, D.; Parrinello, M. *Phys. Rev. Lett.* **2001**, *87*.
- (13) Burger, S. K.; Yang, W. T. *J. Chem. Phys.* **2006**, *124*, 054109.
- (14) Burger, S. K.; Yang, W. T. *J. Chem. Phys.* **2007**, *127*, 164107.
- (15) E, W.; Ren, W.; Vanden-Eijnden, E. *Phys. Rev. B: Condens. Matter Mater. Phys.* **2002**, *66*, 052301.
- (16) E, W.; Ren, W.; Vanden-Eijnden, E. *J. Phys. Chem. B* **2005**, *109*, 6688–6693.
- (17) E, W.; Ren, W.; Vanden-Eijnden, E. *J. Chem. Phys.* **2007**, *126*, 164103.
- (18) Maragliano, L.; Fischer, A.; Vanden-Eijnden, E.; Ciccotti, G. *J. Chem. Phys.* **2006**, *125*, 024106.
- (19) Maragliano, L.; Vanden-Eijnden, E. *Chem. Phys. Lett.* **2007**, *446*, 182–190.
- (20) Khavrutskii, I. V.; Arora, K.; Brooks, C. L. *J. Chem. Phys.* **2006**, *125*, 7.
- (21) Woodcock, H.; Hodoscek, M.; Sherwood, P.; Lee, Y.; Schaefer, H.; Brooks, B. *Theor. Chem. Acc.* **2003**, *109*, 140–148.
- (22) Elber, R. *Curr. Opin. Struct. Biol.* **2005**, *15*, 151–156.
- (23) Karplus, M.; Kuriyan, J. *Proc. Natl. Acad. Sci. U.S.A.* **2005**, *102*, 6679–6685.
- (24) Sheppard, D.; Terrell, R.; Henkelman, G. *J. Chem. Phys.* **2008**, *128*, 10.
- (25) Elber, R.; Ghosh, A.; Cardenas, A. *Acc. Chem. Res.* **2002**, *35*, 396–403.
- (26) Elber, R.; Meller, J.; Olender, R. *J. Phys. Chem. B* **1999**, *103*, 899–911.
- (27) Kabsch, W. *Acta. Crystallogr., Sect. A: Cryst. Phys., Diffraction, Theor. Gen. Crystallogr.* **1976**, *32*, 922–923.
- (28) Allen, M. P.; Tildesley, D. J. *Computer Simulation of Liquids*; Oxford: New York, 1987.
- (29) Press, W. H.; Teukolsky, S. A.; Vetterling, W. T.; Flannery, B. P. *Numerical Recipes in FORTRAN*; 2 ed.; Cambridge University Press: New York, 1992.
- (30) Landau, L. D.; Lifshitz, E. M. *Mechanics*; Elsevier: Oxford, U.K., 1976.
- (31) Evans, D. J.; Morriss, G. P. *Statistical Mechanics of NonEquilibrium Liquids*; Cambridge University Press: New York, 1990.
- (32) Vanden-Eijnden, E.; Heymann, M. *J. Chem. Phys.* **2008**, *128*, 061103.
- (33) Brooks, B. R.; Brucoleri, R. E.; Olafson, B. D.; States, D. J.; Swaminathan, S.; Karplus, M. *J. Comput. Chem.* **1983**, *4*, 187–217.
- (34) Mackerell, A. D.; Bashford, D.; Bellott, M.; Dunbrack, R. L.; Evanseck, J. D.; Field, M. J.; Fischer, S.; Gao, J.; Guo, H.; Ha, S.; Joseph-McCarthy, D.; Kuchnir, L.; Kuczera, K.; Lau, F. T. K.; Mattos, C.; Michnick, S.; Ngo, T.; Nguyen, D. T.; Prodhom, B.; Reiher, W. E.; Roux, B.; Schlenkrich, M.; Smith, J. C.; Stote, R.; Straub, J.; Watanabe, M.; Wiorkiewicz-

- Kuczera, J.; Yin, D.; Karplus, M. *J. Phys. Chem. B* **1998**, *102*, 3586–3616.
- (35) Mackerell, A. D.; Feig, M.; Brooks, C. L. *J. Comput. Chem.* **2004**, *25*, 1400–15.
- (36) Biarnes, X.; Ardevol, A.; Planas, A.; Rovira, C.; Laio, A.; Parrinello, M. *J. Am. Chem. Soc.* **2007**, *129*, 10686–10693.
- (37) Momany, F. A.; Appell, M.; Willett, J. L.; Bosma, W. B. *Carbohydr. Res.* **2005**, *340*, 1638–1655.
- (38) Barrows, S. E.; Dulles, F. J.; Cramer, C. J.; French, A. D.; Truhlar, D. G. *Carbohydr. Res.* **1995**, *276*, 219–251.
- (39) Barrows, S. E.; Storer, J. W.; Cramer, C. J.; French, A. D.; Truhlar, D. G. *J. Comput. Chem.* **1998**, *19*, 1111–1129.
- (40) Guvench, O.; Greene, S. N.; Kamath, G.; Brady, J. W.; Venable, R. M.; Pastor, R. W.; Mackerell, A. D. *J. Comput. Chem.* **2008**, *29*, 2543–64.
- (41) *Eur J Biochem* IUPAC-IUB Joint Commission on Biochemical Nomenclature (JCBN) **1980**, *111*, 295–8.
- (42) Appell, M.; Strati, G.; Willett, J. L.; Momany, F. A. *Carbohydr. Res.* **2004**, *339*, 537–51.
- (43) Balbirnie, M.; Grothe, R.; Eisenberg, D. S. *Proc. Natl. Acad. Sci. U.S.A.* **2001**, *98*, 2375–80.
- (44) Baskakov, I. V.; Legname, G.; Prusiner, S. B.; Cohen, F. E. *J. Biol. Chem.* **2001**, *276*, 19687–90.
- (45) Pan, K. M.; Baldwin, M.; Nguyen, J.; Gasset, M.; Serban, A.; Groth, D.; Mehlhorn, I.; Huang, Z.; Fletterick, R. J.; Cohen, F. E. *Proc. Natl. Acad. Sci. U.S.A.* **1993**, *90*, 10962–6.
- (46) Lipfert, J.; Franklin, J.; Wu, F.; Doniach, S. *J. Mol. Biol.* **2005**, *349*, 648–58.
- (47) Haberthuer, U.; Caflisch, A. *J. Comput. Chem.* **2008**, *29*, 701–715.
- (48) Koslover, E. F.; Wales, D. J. *J. Chem. Phys.* **2007**, *127*.
- (49) Bolhuis, P.; Chandler, D.; Dellago, C.; Geissler, P. *Annu. Rev. Phys. Chem.* **2002**, *53*, 291–318.
- (50) Dellago, C.; Bolhuis, P.; Geissler, P. *Adv. Chem. Phys.* **2002**, *123*, 1–78.
- (51) Dellago, C.; Bolhuis, P. G. *Top. Curr. Chem.* **2007**, *268*, 291–317.

CT9001398