

Surface-Accelerated String Method for Locating Minimum Free Energy Paths

Timothy J. Giese, Şölen Ekesan, Erika McCarthy, Yujun Tao, and Darrin M. York*



Cite This: *J. Chem. Theory Comput.* 2024, 20, 2058–2073



Read Online

ACCESS |



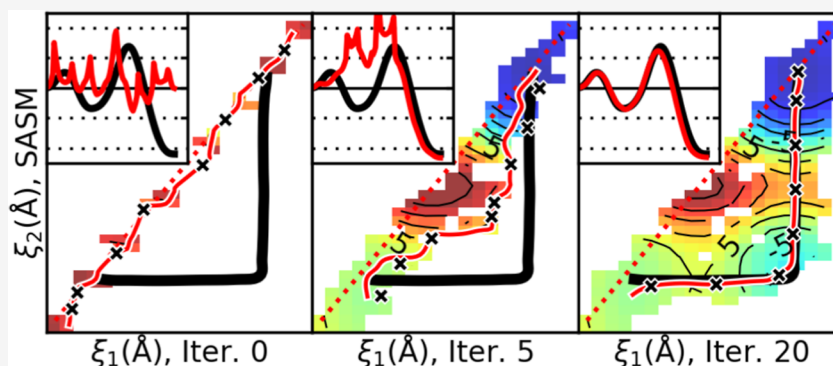
Metrics & More



Article Recommendations



Supporting Information



ABSTRACT: We present a surface-accelerated string method (SASM) to efficiently optimize low-dimensional reaction pathways from the sampling performed with expensive quantum mechanical/molecular mechanical (QM/MM) Hamiltonians. The SASM accelerates the convergence of the path using the aggregate sampling obtained from the current and previous string iterations, whereas approaches like the string method in collective variables (SMCV) or the modified string method in collective variables (MSMCV) update the path only from the sampling obtained from the current iteration. Furthermore, the SASM decouples the number of images used to perform sampling from the number of synthetic images used to represent the path. The path is optimized on the current best estimate of the free energy surface obtained from all available sampling, and the proposed set of new simulations is not restricted to being located along the optimized path. Instead, the umbrella potential placement is chosen to extend the range of the free energy surface and improve the quality of the free energy estimates near the path. In this manner, the SASM is shown to improve the exploration for a minimum free energy pathway in regions where the free energy surface is relatively flat. Furthermore, it improves the quality of the free energy profile when the string is discretized with too few images. We compare the SASM, SMCV, and MSMCV using 3 QM/MM applications: a ribozyme methyltransferase reaction using 2 reaction coordinates, the 2'-O-transphosphorylation reaction of Hammerhead ribozyme using 3 reaction coordinates, and a tautomeric reaction in B-DNA using 5 reaction coordinates. We show that SASM converges the paths using roughly 3 times less sampling than the SMCV and MSMCV methods. All three algorithms have been implemented in the FE-ToolKit package made freely available.

1. INTRODUCTION

The ability to model chemical reactions in the condensed phase¹ using molecular simulations has far-reaching implications for the study of catalysis in biological systems.^{2,3} Advances in fast, accurate quantum mechanical force fields^{4,5} and machine learning models^{6–11} have greatly extended the scope of applications that can be routinely addressed. Nonetheless, simulations of complex reaction pathways remain computationally intensive, and the ongoing development of new methods to improve the robustness and computational cost are important.

Reaction mechanisms can be characterized by calculating a free energy surface in a set of relevant reaction coordinates, the determination of the minimum free energy profile (MFEP) through the surface, and the identification of key stationary points along the MFEP. Many methods for calculating free

energy surfaces have been developed. These approaches can be categorized as¹² methods which analyze equilibrium statistics obtained from umbrella sampling,^{13–16} methods which analyze nonequilibrium statistics^{17–20} based on the work of Jarzynski,²¹ and methods that integrate auxiliary degrees of freedom, such as λ -dynamics^{22–25} and metadynamics.^{26,27} Similarly, there are two general approaches for locating a minimum free energy path.²⁸ The first approach is to sample the reaction over

Received: December 22, 2023

Revised: February 6, 2024

Accepted: February 7, 2024

Published: February 17, 2024



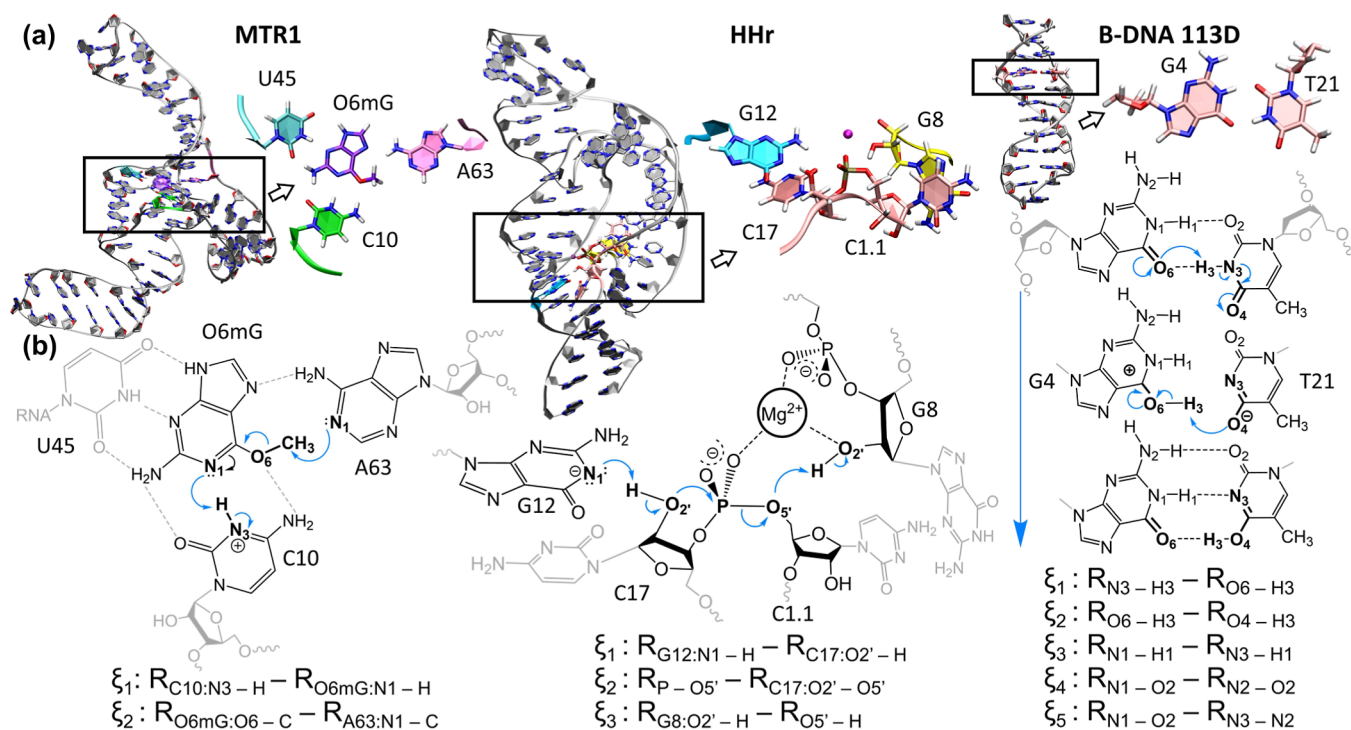


Figure 1. (a) MTR1 ribozyme, HHr ribozyme, and a B-DNA with a GT wobble pair examined in this work. The rectangles highlight the active site region. (b) Reaction mechanisms and reaction coordinates. The B-DNA system is a tautomer reaction which transfers the T21 N3 proton to position O4 and reorganization of the G:T hydrogen bond network. The shown atomic configurations correspond to the reactant state. The black and gray atoms denote the QM region and nearby MM atoms, respectively.

a wide range of reaction coordinate values to obtain a relatively complete picture of the free energy surface through which a path can be optimized. This approach can greatly benefit from enhanced sampling methods, such as replica exchange molecular dynamics²⁹ and integrated tempering sampling.³⁰ The second, and more cost-effective, approach is to use a chain-of-states method, such as nudged elastic band³¹ or the string method,^{32,33} to direct the sampling toward the MFEP, thereby reducing the amount of effort spent simulating irrelevant, high-energy regions of the free energy surface.

Many variations of the string method^{34–41} have been developed that are capable of being applied to large-scale problems, like protein folding.^{36,37} These applications often describe the path using a large number of reaction coordinates,⁴² direct comparison of Cartesian coordinates,³² path collective variables,^{28,43} the use of the hills method,^{26,44} or machine learning techniques.⁴⁵ Although string method development was originally motivated by the desire to use many reaction coordinates,^{34,35,38} many examples can be found of their use in quantum mechanical/molecular mechanical (QM/MM) applications involving only a few reaction coordinates.^{46–52} String methods, such as the one presented in ref 38, are particularly appealing because it is performed with standard umbrella sampling with harmonic biasing potentials, which are widely supported across simulation packages. Because QM/MM sampling is very costly, the present work seeks to optimize the string method described in ref 38 specifically for cases involving QM/MM simulations with a few reaction coordinates. The new method reduces the number of string iterations required to reach convergence because it uses the current estimate of the unbiased free energy to accelerate the exploration of flat regions of the surface. In this respect, the new method draws inspiration from ideas

behind the metadynamics approach,^{26,27} however, the new method only requires sampling obtained using standard harmonic biasing potentials.

We describe a new surface-accelerated string method (SASM) and compare it to two similar algorithms: the string method in collective variables^{34,35} (SMCV), and the modified string method in collective variables³⁸ (MSMCV). We have implemented all 3 of these methods in the *ndfes* software⁴⁹ freely distributed within the FE-ToolKit package.⁵³ The FE-ToolKit package has also been incorporated in the open-source AmberTools simulation suite.⁵⁴ There are several key differences between the SASM and related string methods. First, the SASM is a hybrid of the two approaches for locating an MFEP (chain-of-states method versus calculation of a multidimensional free energy surface). Whereas the SMCV and MSMCV update the path from the sampling obtained in the most recent string iteration, the SASM optimizes the path on the current estimate of the multidimensional free energy surface calculated from the aggregate sampling of all string iterations. In this respect, the SASM is similar to some adaptive umbrella sampling strategies.⁵⁵ Second, the SASM decouples the number of images used to represent the path from the number of simulated images. The SMCV and MSMCV methods construct a new path by fitting a curve that interpolates a set of discrete control points obtained from a corresponding number of simulated images; therefore, if there was an insufficient number of images, the path may cut corners. By decoupling the representation of the path from the number of simulated images, the level of detail used to describe the path is not limited by the number of simulations. Third, unlike the SMCV and MSMCV, the SASM does not require the images to be simulated along the current estimate of the path. We take advantage of this by introducing alternating stages of

“exploration” and “refinement” steps. The exploration steps propose new simulations offset from the path in the direction that the path is moving, and the refinement steps place simulations along the path in a manner that improves the phase space overlap.

We compare the progress of the string optimizations using the SMCV, MSMCV, and SASM with respect to the number of simulations per string, the sampling per simulation, and the spline representation of the path (either piecewise linear or Akima spline paths) in 3 applications. The first application uses 2 reaction coordinates to describe a ribozyme undergoing a methyl transfer reaction (MTR1)^{56–58} (PDB ID 7V9E). The second application uses 3 reaction coordinates to model the 2'-O-transphosphorylation reaction of Hammerhead ribozyme (HHr)⁵⁹ (PDB ID 2OEU). The third application uses 5 reaction coordinates to optimize a tautomeric reaction pathway in B-DNA (PDB ID 113D).⁶⁰ Schematics of the 3 systems are shown in Figure 1. We demonstrate that the SASM converges the MFEP faster than the SMCV and MSMCV when we vary the amount of sampling. The SASM avoids artifacts that can occur in the path “reparametrization step” of the SMCV and MSMCV. Finally, we show that the SASM method will sample the path in an efficient manner that achieves good overlap between the biased simulations when the number of simulations is reduced.

2. METHODS

2.1. String Method in Collective Variables. This section summarizes the SMCV method, which was originally described in refs 34 and 35. Let \mathbf{x} and $\mathbf{q}(\mathbf{x})$ be the 3N array of atomic positions and N_{dim} reaction coordinate values, respectively. Umbrella sampling is performed at N_{img} images along the path using a biased potential energy function, U_n .

$$U_n(\mathbf{x}) = U(\mathbf{x}) + W(\mathbf{q}(\mathbf{x}); \mathbf{k}_n, \mathbf{q}_n) \quad (1)$$

Image n is biased by a potential W that is parametrized by N_{dim} harmonic force constants \mathbf{k}_n and equilibrium positions \mathbf{q}_n . In other words, N_{dim} is the size of the reduced dimensional space of reaction coordinates.

$$W(\mathbf{q}(\mathbf{x}); \mathbf{k}_n, \mathbf{q}_n) = \frac{1}{2} \sum_{d=1}^{N_{\text{dim}}} k_{nd} (q_d(\mathbf{x}) - q_{nd})^2 \quad (2)$$

The algorithm for calculating the SMCV consists of the following steps:

1. Sample each of the N_{img} images along the path for some amount of time, Δt . The images differ by their biasing potentials, which center the harmonic potentials at discrete points along the current estimate of the path, \mathbf{q}_n .
2. Analyze the sampling to update (evolve) the reaction coordinate values, $\mathbf{q}_{c,n}$. The “control points”, $\mathbf{q}_{c,n}$ are discrete estimates along the new path, but they do not necessarily uniformly discretize it. The calculation of the control points is sometimes called the “evolution step”.
3. Construct a parametric curve that interpolates the control points. The parametric curve is the new estimate of the path.
4. Uniformly discretize the parametric curve to obtain the biasing potential centers for the next iteration. The construction of a new curve and its discretization is sometimes called the “reparametrization step”.

The SMCV evolution step is given by eq 3, where $q_{nd}^{(k)}$ is the value of the reaction coordinate d of image n at string iteration k , and $q_{c,nd}^{(k+1)}$ is a control point used to define the parametric curve in string iteration $k + 1$, discussed in the next section. Each image is simulated for a length of time Δt , and $\langle \cdot \rangle_{\mathbf{k}_n, \mathbf{q}_n}$ denotes a time average obtained from image n .

$$q_{c,nd}^{(k+1)} = q_{nd}^{(k)} - \frac{\Delta t}{\gamma} \sum_{d'=1}^{N_{\text{dim}}} M_{dd'}(\mathbf{k}_n, \mathbf{q}_n^{(k)}) \nabla G_{nd'}(\mathbf{k}_n^{(k)}, \mathbf{q}_n^{(k)}) \quad (3)$$

$\nabla G_{nd}(\mathbf{k}_n^{(k)}, \mathbf{q}_n^{(k)})$ approximates the free energy gradient about the point $\mathbf{q}_n^{(k)}$ in dimension d .

$$\nabla G_{nd}(\mathbf{k}_n, \mathbf{q}_n) = - \left\langle \frac{\partial W(\mathbf{q}(\mathbf{x}); \mathbf{k}_n, \mathbf{q}_n)}{\partial q_{nd}} \right\rangle_{\mathbf{k}_n, \mathbf{q}_n} \quad (4)$$

\mathbf{M} is closely related to a product of mass-weighted Wilson B-matrices;⁶¹ that is to say, $\nabla_a q$ is the gradient of the reaction coordinate value with respect to the atomic positions of the atom a and m_a is an atomic mass.

$$M_{dd'}(\mathbf{k}_n, \mathbf{q}_n) = \left\langle \sum_{a=1}^N \frac{\nabla_a q_d(\mathbf{x}) \cdot \nabla_a q_{d'}(\mathbf{x})}{m_a} \right\rangle_{\mathbf{k}_n, \mathbf{q}_n} \quad (5)$$

γ is a friction coefficient, a parameter of the method. The numerical stability of the SMCV critically depends on the ratio $\Delta t \gamma^{-1}$. In ref 35, it was found that the method was stable when choosing $\gamma = 1500 \text{ ps}^{-1}$ when $\Delta t = 20 \text{ fs}$. In the present work, we adjust γ to maintain this same ratio when Δt is varied. The construction of parametric curves and their uniform discretization are described in the next section.

2.2. Parametric Curves and the Reparametrization Step. We represent a continuous path as a parametric curve of reaction coordinates, $\mathbf{q}(p)$, where $p \in [0, 1]$ is a progress variable such that $p = 0$ and $p = 1$ denote two ends of the path. In other words, the path at string iteration k , $\mathbf{q}^{(k)}(p)$ is an array of N_{dim} one-dimensional splines that are chosen such that each spline interpolates the N_{img} control points, $q_{c,nd}^{(k)}$ located at a common set of progress control values, $p_{c,n}^{(k)}$. Let $\mathbf{p}_c^{(k)} = \{p_{c,1}^{(k)}, \dots, p_{c,N_{\text{img}}}^{(k)}\}$ and $\mathbf{q}_{c,d}^{(k)} = \{q_{c,1,d}^{(k)}, \dots, q_{c,N_{\text{img},d}}^{(k)}\}$ denote the $N_{\text{img}} \times 1$ arrays of progress control values and control points in dimension d , respectively. The spline representation of the path in dimension d , $q_d^{(k)}(p)$ is parametrized from these quantities.

$$q_d^{(k)}(p) \equiv q_d(p; \mathbf{q}_{c,d}^{(k)}, \mathbf{p}_c^{(k)}) \quad (6)$$

In the context of the SMCV (or similar string methods), the control points are the new estimates of the reaction coordinates after the evolution step (eq 3). In some cases, one may choose to reduce the numerical noise in the path by first applying a smoothing procedure, in which case the control points are the reaction coordinate values after smoothing. The results presented in this work use a smoothing algorithm implemented in the `ndfes` software when the parametric curve is modeled with Akima spline functions,⁶² but we do not apply smoothing to the control points when using piecewise linear paths. The details of the smoothing algorithm are described in the Supporting Information.

The parametric curve depends on the progress control values, which are interpreted as fractional arc lengths through the curve. If the path is a piecewise linear function connecting

the control points, then the progress control values can be calculated from the Euclidean distance between adjacent points, as shown in eq 7.

$$p_{c,n}^{(k)} \approx \begin{cases} 0, & \text{if } n = 1 \\ \frac{\sum_{m=2}^n \sqrt{\sum_{d=1}^{N_{\text{dim}}} (q_{md}^{(k)} - q_{m-1,d}^{(k)})^2}}{\sum_{m=2}^{N_{\text{img}}} \sqrt{\sum_{d=1}^{N_{\text{dim}}} (q_{md}^{(k)} - q_{m-1,d}^{(k)})^2}}, & \text{otherwise} \end{cases} \quad (7)$$

Alternatively, if the parametric curve is a set of Akima spline functions⁶² (or any smooth interpolating function), then eq 7 is only an approximation of the fractional arc lengths. Accurate values of the progress control values can be found by iteratively solving eq 8.

$$p_{c,n}^{(k,i+1)} = \frac{\int_0^{p_{c,n}^{(k,i)}} \sqrt{\sum_{d=1}^{N_{\text{dim}}} \left(\frac{\partial q_d(p; \mathbf{q}_{c,d}^{(k,i)}, \mathbf{p}_c^{(k,i)})}{\partial p} \right)^2} dp}{\int_0^1 \sqrt{\sum_{d=1}^{N_{\text{dim}}} \left(\frac{\partial q_d(p; \mathbf{q}_{c,d}^{(k,i)}, \mathbf{p}_c^{(k,i)})}{\partial p} \right)^2} dp} \quad (8)$$

eq 8 describes the iterative solution of $p_{c,n}^{(k)}$ by introducing a second, auxiliary index i , $p_{c,n}^{(k,i)}$. The initial values, $p_{c,n}^{(k,0)}$ are the piecewise linear approximation shown in eq 7, and we terminate the iterative solution when $\sum_{m=1}^{N_{\text{img}}} (p_{c,m}^{(k,i+1)} - p_{c,m}^{(k,i)})^2 < 10^{-16}$. We drop the auxiliary index to denote the converged progress control values.

Given the parametric spline representation of the path, the uniformly discretized images for string iteration $k + 1$ are shown in eq 9, where $p_n = (n - 1)/(N_{\text{img}} - 1)$.

$$q_{nd}^{(k+1)} = q_d(p_n; \mathbf{q}_{c,d}^{(k+1)}, \mathbf{p}_c^{(k+1)}) \quad (9)$$

2.3. Modified String Method in Collective Variables.

The MSMCV was originally presented in ref 38; it differs from the SMCV only by replacing the evolution step (eq 3) with eq 10.

$$q_{c,nd}^{(k+1)} = \langle q_d(\mathbf{x}) \rangle_{\mathbf{k}_n^{(k)}, \mathbf{q}_n^{(k)}} \quad (10)$$

In other words, the control points for the new path are the mean observed positions of the reaction coordinates from the simulations performed along the current path. Upon finding the control points, a new parametric curve is fit. The curve is uniformly discretized to define the new positions of the biasing potentials.

2.4. Surface-Accelerated String Method. The SASM constructs a N_{dim} dimensional free energy surface from the available sampling and optimizes a path on that surface. A decision is then made to place a new set of simulations, which may or may not be along the optimized path. When the new simulations are placed along the path, we refer to it as a “refinement step”. Alternatively, we allow for “exploration steps” that offset the simulations from the path in the direction that the path is moving.

The algorithm for calculating the SASM consists of the following steps.

1. Sample each of the N_{img} images for some amount of time.
2. Construct a N_{dim} dimensional unbiased free energy surface by analyzing the aggregate sampling produced from all simulations and string iterations. This is the best

estimate of the free energy surface from the available sampling. The N_{dim} dimensional space is discretized into bins, and the free energy value and the number of observed samples in each bin are tabulated.

3. Create a smooth representation of the free energy surface, such that the free energy value and gradient can be readily computed at any point in the space of reaction coordinates.
4. Use the free energy surface to optimize an MFEP in the space of reaction coordinates. This optimization procedure does not involve the generation of additional sampling. Instead, the optimization is performed on a fixed free energy surface using a series of “synthetic string iterations”, described below.
5. If the current iteration is an even integer, then place the new simulations along the path. If the current iteration is an odd integer, then allow the new set of simulations to be displaced from the path by some amount in the direction that the path is moving.

The unbiased free energy can be calculated using established methods, such as the variational free energy profile method,^{49,63,64} the multistate Bennett acceptance ratio (MBAR) method,⁶⁵ or the unbinned weighted histogram (UWHAM) method.^{66,67} As discussed in ref 49, a smooth representation of the free energy surface can be made using one of many methods, including the use of Cardinal B-Splines,⁶⁸ radial basis functions,^{69,70} or Gaussian process regression.⁷¹ In the present work, we calculate the free energy surface by solving the MBAR/UWHAM equations to reweight the biased sampling. The samples are histogrammed, and the free energy of each bin is tabulated. We use fourth-order Cardinal B-splines to represent the surface as a smooth function. A mathematical description of the B-spline interpolation is provided in the Supporting Information for completeness. The free energy values are formally defined only in those regions whose histogram bins are occupied by at least one sample. In practice, we exclude all bins containing fewer than 10 samples because their free energy values are often unreliable.

To optimize a path on a fixed free energy surface, we adapt the MSMCV by replacing eqs 10 with 11, where $F(\mathbf{q})$ is the value of the unbiased free energy at \mathbf{q} .

$$\mathbf{q}_{c,n}^{(k+1,s+1)} = \arg \min_{\mathbf{q}} \{F(\mathbf{q}) + W(\mathbf{q}; \bar{\mathbf{k}}^{(k)}, \mathbf{q}_n^{(k+1,s)})\} \quad (11)$$

$$\mathbf{q}_n^{(k+1,s+1)} = \mathbf{q}(p_n; \mathbf{q}_c^{(k+1,s+1)}, \mathbf{p}_c^{(k+1,s+1)}) \quad (12)$$

The $q_{c,nd}^{(k,s)}$ values are the control points of the synthetic images used to describe the path. Specifically, n indexes the synthetic image, d indexes the dimension, k is the string iteration, and s is the synthetic iteration. The number of synthetic images, N_{simgr} does not need to be the same as the number of images used to perform explicit simulations, N_{img} . In the present work, we use $N_{\text{simgr}} = 100$ to describe the path. The biasing potential appearing in eq 11 requires a set of force constants for each synthetic image. If the number of simulated images was the same as the number of synthetic images, then the simulation force constants could be reused to define the synthetic image biasing potentials. The number of synthetic images is often much larger than the number of simulated images; therefore, one needs to transform the $N_{\text{img}} \times N_{\text{dim}}$ simulation force constants into a set of $N_{\text{simgr}} \times N_{\text{dim}}$ force constants. Our choice

is to use a common set of force constants for each synthetic image by averaging the simulated force constants; that is, $\bar{k}_d = N_{\text{img}}^{-1} \sum_{n=1}^{N_{\text{img}}} k_{nd}^{(k)}$. Equation 11 is analogous to the MSMCV, but instead of performing a biased simulation of $3N$ atomic coordinates to obtain the reaction coordinate distribution means, one performs a minimization directly on a biased N_{dim} free energy surface. In other words, eq 11 is a synthetic iteration that allows us to repeatedly propagate the string without producing additional sampling. The path is optimized with N_{siter} iterations (or until convergence is sufficiently met), such that $\mathbf{q}_{\text{opt}}^{(k+1)}(p) \equiv \mathbf{q}(p; \mathbf{q}_c^{(k+1, N_{\text{siter}})}, \mathbf{p}_c^{(k+1, N_{\text{siter}})})$ is the best estimate of the MFEP from the available sampling. The optimized synthetic control points also serve as the initial guess for the path in the next string iteration: $\mathbf{q}_{c,n}^{(k+1,0)} = \mathbf{q}_{c,n}^{(k, N_{\text{siter}})}$.

By optimizing the MFEP on the current estimate of the free energy surface, the N_{img} real images are no longer responsible for describing the path. Instead, their sole responsibility is to provide sampling to improve the quality and range of the free energy surface. For this purpose, the SASM evolution step (eq 13) includes two modifications relative to a simple uniform discretization.

$$\mathbf{q}_n^{(k+1)} = \mathbf{q}_{\text{opt}}^{(k+1)}(p_n + \Delta p^{(k+1)}) + \Delta \mathbf{q}_n^{(k+1)} \quad (13)$$

The first modification is a shifting of the progress control points when discretizing the parametric curve

$$p_n + \Delta p^{(k+1)} = \begin{cases} p_{n,0}, & \text{if } N(p_{n,0}) = 0 \\ p_{n,-1/3}, & \text{else if } N(p_{n,-1/3}) = 0 \\ p_{n,+1/3}, & \text{else if } N(p_{n,+1/3}) = 0 \\ p_{n,0}, & \text{else if } \text{mod}(k, 3) = 0 \\ p_{n,-1/3}, & \text{else if } \text{mod}(k, 3) = 1 \\ p_{n,+1/3}, & \text{else if } \text{mod}(k, 3) = 2 \end{cases} \quad (14)$$

where $p_{n,0}$ is a uniform discretization, and $p_{n,+1/3}$ and $p_{n,-1/3}$ shift the discretization by 1/3 of the distance to a neighboring image.

$$p_{n,x} = \max \left(0, \min \left(1, \frac{n-1+x}{N_{\text{img}}-1} \right) \right) \quad (15)$$

$N(p)$ is the number of samples that have been observed at the point $\mathbf{q}_{\text{opt}}^{(k+1)}(p)$. In other words, the first 3 cases in eq 14 check whether there are gaps in the sampling along the path. If there is a gap, then sampling at that position is prioritized. The last 3 cases in eq 14 are a schedule that is followed when no gaps in the sampling are detected. The schedule alternates between these displacements during the course of the string optimization to help ensure that one obtains sufficient sampling along the path in the event that one underestimates an appropriate value of N_{img} .

The second modification is the introduction of $\Delta \mathbf{q}_n^{(k+1)}$ which displaces the image in the direction of the path's movement.

$$\Delta \mathbf{q}_n^{(k+1)} = \begin{cases} 0, & \text{if } \text{mod}(k, 4) = 0 \\ [1 - \delta_{0, N(p_n + \Delta p^{(k+1)})}] \frac{\Delta \mathbf{q}_{\text{min},n}^{(k+1)}}{|\Delta \mathbf{q}_{\text{min},n}^{(k+1)}|} h_1, & \text{if } \text{mod}(k, 4) = 1 \\ 0, & \text{if } \text{mod}(k, 4) = 2 \\ [1 - \delta_{0, N(p_n + \Delta p^{(k+1)})}] h_2 \frac{\Delta \mathbf{q}_{\text{min},n}^{(k+1)}}{|\Delta \mathbf{q}_{\text{min},n}^{(k+1)}|}, & \text{if } \text{mod}(k, 4) = 3 \end{cases} \quad (16)$$

We refer to $\Delta \mathbf{q}_n^{(k+1)} = 0$ as a refinement step that places the simulations along the path, and the other cases are exploration steps intended to better describe the free energy surface in the vicinity of the path in the direction of its movement. The exploration steps accelerate the evolution of the string through flat areas of the free energy surface. The leading Kronecker delta function causes the exploration step to be skipped if a gap in the sampling was previously detected in eq 14. The exploration direction is determined from the difference between the optimized paths of the current and previous iterations.

$$\Delta \mathbf{q}_{\text{min},n}^{(k+1)} = \mathbf{q}_{\text{opt}}^{(k+1)}(p_n + \Delta p^{(k+1)}) - \mathbf{q}_{\text{opt}}^{(k)}(p^*) \quad (17)$$

The value of p^* is the point on the previous path that is closest to the point $p_n + \Delta p^{(k+1)}$ on the current path.

$$p^* = \arg \min_p |\mathbf{q}_{\text{opt}}^{(k+1)}(p_n + \Delta p^{(k+1)}) - \mathbf{q}_{\text{opt}}^{(k)}(p)|^2 \quad (18)$$

The h_m values are the magnitude of the displacement, where w_d is the width assigned to each dimension. In the present work, we use $w_d = 0.15 \text{ \AA}$ for all dimensions, which is also the width of the histogram bins used to construct the free energy surface.

$$h_m = \min \left(mw_1 \frac{|\Delta \mathbf{q}_{\text{min},n}^{(k+1)}|}{|\Delta \mathbf{q}_{\text{min},n,1}^{(k+1)}|}, \dots, mw_{N_{\text{dim}}} \frac{|\Delta \mathbf{q}_{\text{min},n}^{(k+1)}|}{|\Delta \mathbf{q}_{\text{min},n,N_{\text{dim}}}^{(k+1)}|} \right) \quad (19)$$

If one imagines the point at $p_n + \Delta p^{(k+1)}$ as being located at the corner of a voxel, then eq 19 can be interpreted as choosing the magnitude to be the maximum displacement that does not exceed the range of m voxels.

The SASM method has several parameters that can be adjusted, including the number of simulated images N_{img} , the number of synthetic images N_{sim} , the simulation length Δt , the biasing potential force constants \mathbf{k} , the cyclic schedules shown in eqs 14 and 16, and the voxel width w_d appearing in eq 19. Unlike the SMCV and MSMCV, the SASM requires an estimate of the unbiased free energy surface to propagate the string. The calculation of an unbiased free energy surface from a solution of the MBAR/UWHAM equations requires overlap between the biased distributions,^{65–67} so a suitably large value of N_{img} is necessary. An optimal choice of N_{img} depends on the length of the string, the gradient of the underlying free energy surface, and the values of \mathbf{k} and Δt . In practice, one can validate their choice of N_{img} by analyzing the sampling overlap produced from their initial guess pathway. In the event that the

MFEP was significantly longer than the initial guess, N_{img} may become too small during the course of the string method. We did not design the SASM to dynamically choose the value of N_{img} based on the current string length because changes in N_{img} may complicate resource allocation scheduling requests. Instead, the progress variable shifts (eq 14) effectively increase N_{img} over the course of several string iterations to minimize the consequences of having chosen an insufficiently small value. Using a cyclic schedule of length 3, eq 14 acts to effectively increase N_{img} by a factor of 3 over a series of string iterations. Alternatively, if the calculations were prepared with an excessive number of images, then the proposed shifts would be very small in relation to the length of the string, $\mathbf{q}_{\text{opt}}^{(k+1)}(p_n + \Delta p^{(k+1)}) \approx \mathbf{q}_{\text{opt}}^{(k+1)}(p_n)$, rendering the shifts unnecessary.

The displacement schedule (eq 16) alternates between refinement and exploration iterations to avoid introducing a bias to the exploration direction (eq 17). In other words, if the exploration direction was chosen by optimizing a path on a surface that introduced new sampling in areas offset from the current path, then the new path is more likely to move in the direction of the added sampling. Our recommended displacement schedule cycles every 4 iterations rather than 2 iterations. We have explored the use of 2 iteration schedules that alternate between refinement and exploration iterations with h_1 or h_2 (eq 19) using $w_d = 0.15 \text{ \AA}$. In brief, the 2 iteration schedule involving h_2 was found to converge the MFEP faster than the 2 iteration schedule involving h_1 , and it performed about as well as the 4 iteration schedule. We recommend the 4 iteration schedule because it is less likely to produce gaps in the sampling. Appropriate values of h_m are coupled to the value of w_d . We use histogram bin widths of 0.15 \AA to construct the free energy surface because smaller bin widths are more likely to produce numerical noise in the free energy surface, which manifests as noise in the optimized path.⁴⁹ Furthermore, if the displacements become too large, then extended simulations may be necessary to equilibrate the system after making a significant change to the biasing potential. Once the SASM has converged, further iterations will fluctuate around the MFEP. Consequently, the sampling will envelop the MFEP.

2.5. Computational Details. All QM/MM simulations in this work were performed with the sander molecular dynamics software⁵⁴ using the default leapfrog integrator with a 1 fs integration time step. More efficient sampling with a longer integration time step may be obtained using the recently developed “middle” thermostat scheme described in refs 72 and 73., which is already available in the Amber software. The SHAKE algorithm⁷⁴ was used to fix MM bonds involving hydrogen, whereas all QM bonds were left unconstrained. The covalent bonds at the QM/MM boundary were capped with the hydrogen link-atom approach.^{75,76} Electrostatics were calculated with the particle mesh Ewald method^{77–79} adapted for use within semiempirical QM/MM simulations^{80,81} using tinfoil boundary conditions^{82,83} a 1 \AA^3 reciprocal space grid, and 10 \AA real space cutoffs. The Lennard-Jones interactions were similarly calculated to 10 \AA and a long-range tail correction was included to account for the interactions beyond the cutoff.⁸⁴

The MTR1 ribozyme (PDB ID 7V9E⁵⁸) consists of 2207 atoms with a net 66–charge. The ribozyme was solvated with 18,250 TIP4P/Ew waters, 113 sodium ions, and 47 chlorine ions in a truncated octahedron with real space lattice vectors of length 90.2 \AA resulting in 75,367 particles and an ion

concentration of 140 mM. The ff99OL3 RNA force field⁸⁵ and Joung and Cheatham⁸⁶ monovalent ion parameters have been used. Details regarding the preparation and equilibration of this system have already been reported elsewhere.⁸⁷ In brief, the pressure and temperature were equilibrated for 50 ns with the MM force field potential to maintain 1 atm and 298 K in the isothermal–isobaric ensemble using the Berendsen barostat⁸⁸ and Langevin thermostat⁸⁹ with a collision frequency of 5 ps^{-1} . At this point, the MM force field was replaced with the DFTB3 QM/MM potential using the “3ob” parameter set.⁹⁰ The QM region consists of 48 atoms with net 1+ charge, as illustrated in Figure 1. A QM/MM simulation of the reactant state was equilibrated for 12.5 ps in the canonical ensemble at 298 K. The DFTB3 QM/MM umbrella production sampling was similarly performed at constant temperature with $200 \text{ kcal mol}^{-1} \text{ \AA}^{-2}$ force constants on the two reaction coordinates describing the transfer of a proton, $\xi_1 = R_{\text{C10:N3-H}} - R_{\text{O6mG:N1-H}}$ and methyl group, $\xi_2 = R_{\text{O6mG:O6-C}} - R_{\text{A63:N1-C}}$ as visualized in Figure 1.

The HHr ribozyme (PDB ID 2OEU⁵⁹) consists of 2020 atoms with a net 62–charge. The ribozyme was solvated with 13,319 TIP4P/Ew waters, 5 magnesium ions (replacing the crystal structure manganese ions), 86 sodium ions, and 34 chlorine ions in a truncated octahedron with real space lattice vectors of length 81.7 \AA resulting in 55,421 particles and an ion concentration of 140 mM. The ff99OL3 RNA force field,⁸⁵ Joung and Cheatham⁸⁶ monovalent ion, and Li-Merz⁹¹ 12-6-4 divalent ion parameters with Panteva^{92,93} corrections, which ensure balanced interactions between metal ions and nucleic acids, have been used. Full details of the preparation and equilibration of this system have been reported elsewhere.⁹⁴ In brief, the pressure and temperature were equilibrated for 100 ns with the MM force field potential to maintain 1 atm and 298 K in the isothermal–isobaric ensemble. The MM force field was replaced with the AM1/d QM/MM potential.⁹⁵ The QM region consists of 85 atoms with net 1–charge. The QM region is illustrated in Figure 1; for clarity, the Mg^{2+} and the 4 waters directly coordinating the Mg^{2+} were included in the QM region. A QM/MM simulation of the reactant state was equilibrated for 50 ps in the canonical ensemble at 298 K. All AM1/d QM/MM umbrella production sampling was performed at constant temperature with $200 \text{ kcal mol}^{-1} \text{ \AA}^{-2}$ force constants on the three reaction coordinates describing the proton transfer from the nucleophile to the general base, $\xi_1 = R_{\text{G12:N1-H}} - R_{\text{C17:O2'-H}}$, and phosphoryl transfer, $\xi_2 = R_{\text{P-O5'}} - R_{\text{C17:O2'-O5'}}$, and the proton transfer from the general acid to the leaving group, $\xi_3 = R_{\text{G8:O2'-H}} - R_{\text{O5'-H}}$.

The B-DNA sequence (PDB ID 113D)⁶⁰ consists of 762 atoms with a net 22–charge. The ribozyme was solvated with 5151 TIP4P/Ew waters, 35 sodium ions, and 13 chlorine ions in a truncated octahedron with real space lattice vectors of length 59.3 \AA resulting in 21,414 particles and an ion concentration of 140 mM. The system was modeled with the OLS DNA force field⁹⁶ and Joung and Cheatham⁸⁶ monovalent ion parameter set. The system was prepared by minimizing the solvent environment and hydrogen positions while restraining the DNA heavy atoms, followed by a gradual heating of the system from 0 to 298 K over the course of 300 ps in the *NVT* ensemble, and the system density was equilibrated at 1 atm for 8 ns in the *NPT* ensemble. The MM force field was replaced with the AM1/d QM/MM potential,⁹⁵ where the QM region (the G4 and T21 nucleobases depicted in Figure 1) consists of 31 atoms with

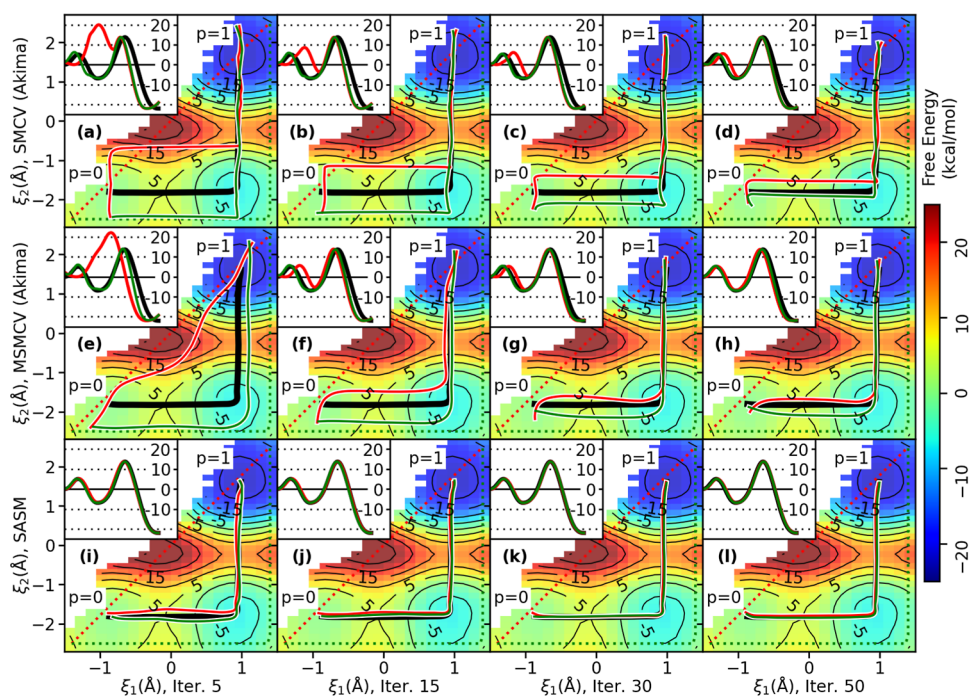


Figure 2. Progress of the string methods at several iterations of the MTR1 reaction starting from concerted (red lines) and stepwise (green lines) initial guess paths. Parts (a–d), (e–h), and (i–l) illustrate the convergence of the SMCV, MSMCV, and SASM, respectively. Each string is composed of 32 images, and each image is sampled for 4 ps. The initial guesses are dashed lines. The colored areas are the best estimate of the free energy surface, calculated from the aggregate sampling produced by all string methods. The black line is the MFEP optimized on the best estimate of the surface. The insets are the reference free energy values along the paths (kcal/mol).

a net neutral charge. A QM/MM simulation of the reactant state was equilibrated for 50 ps in the canonical ensemble at 298 K. All AM1/d QM/MM umbrella production sampling was performed at constant temperature with $200 \text{ kcal mol}^{-1} \text{ \AA}^{-2}$ force constants for each of the reaction coordinates listed in Figure 1.

To start any string method, one must first construct a series of structures to be used as the initial guess. For the MTR1 reaction, we consider two initial guesses: a concerted guess that uniformly discretizes a line connecting the approximate position of the reactant state $\xi_{\text{react.}} = (-1.4, -2.5 \text{ \AA})$ to the product state $\xi_{\text{prod.}} = (1.4, 2.5 \text{ \AA})$, and a stepwise guess that uniformly discretizes a piecewise linear path connecting the reactant state, approximate intermediate state $\xi_{\text{inter.}} = (1.4, -2.5 \text{ \AA})$, and the product state. For the HHR reaction, the initial guess discretizes a linear transformation between the approximate reactant state $\xi_{\text{react.}} = (-1, -2, -1 \text{ \AA})$ to the approximate product state $\xi_{\text{prod.}} = (1, 2, 1 \text{ \AA})$. Similarly, the initial guess for the B-DNA tautomer reaction discretizes a linear transformation between the reactant $\xi_{\text{react.}} = (-0.86, -0.60, -2.0, -0.76, -2.6 \text{ \AA})$ and product states $\xi_{\text{react.}} = (0.67, 0.78, -0.82, 0.84, 0.22 \text{ \AA})$. The atomic coordinates were generated from a sequence of short (200 fs) simulations that restart each image from the final structure of the previous image. After this scan was completed, each image was independently equilibrated for an additional 4 ps. The final coordinates from these equilibrations became the starting structures to initiate the string method.

The SMCV, MSMCV, and SASM were performed multiple times while varying the number of images and length of production sampling. The MTR1 simulations performed for 4 ps/image and 500 fs/image saved 400 samples/image and 250 samples/image, respectively. The HHR simulations performed

for 625 fs/image and 312 fs/image saved 125 samples/image and 156 samples/image, respectively. The B-DNA simulations were performed for 1 ps/image and 200 samples/image were saved. In all cases, we analyze only the last 75% of saved samples when solving the MBAR/UWHAM equations.

3. RESULTS AND DISCUSSION

Here, we compare the SMCV, MSMCV, and SASM string methods using three reactive chemical systems having varying number of reaction coordinates. 1. A 2D example of an artificially engineered methyltransferase ribozyme (MTR1)⁵⁶ that catalyzes the methylation of a target adenine. 2. A 3D example of a naturally occurring HHR⁵⁹ that catalyzes site-specific RNA self-cleavage. 3. A 5D example of tautomerization in dG-dT wobble pairs that lead to misincorporation during replication.⁹⁷

3.1. MTR1 Catalytic Mechanism. Evolutionary theories based on an RNA world^{98,99} presumably would require RNA molecules to catalyze C–C and C–N bond formation essential for nucleic acid synthesis and early metabolic transformations. There are no known naturally occurring examples of RNA enzymes that have this ability. Recently, an MTR1 has evolved *in vitro*⁵⁶ that binds *O*⁶-methylguanine and catalyzes the methylation of a target adenine (A63) at the N1 position^{57,100} (Figure 1a). Computational enzymology studies performed by our group,⁸⁷ in collaboration with Huang, Lilley, and co-workers,⁵⁸ revealed a surprisingly sophisticated mechanism that involves a protonated cytosine residue that acts as an acid to facilitate site-specific C–N bond formation, broadening the range of known RNA-catalyzed chemistry and further demonstrating the versatility of RNA catalysis.¹⁰¹ In the computational study, we employed an early version of the string method and found it to be slowly convergent, making it

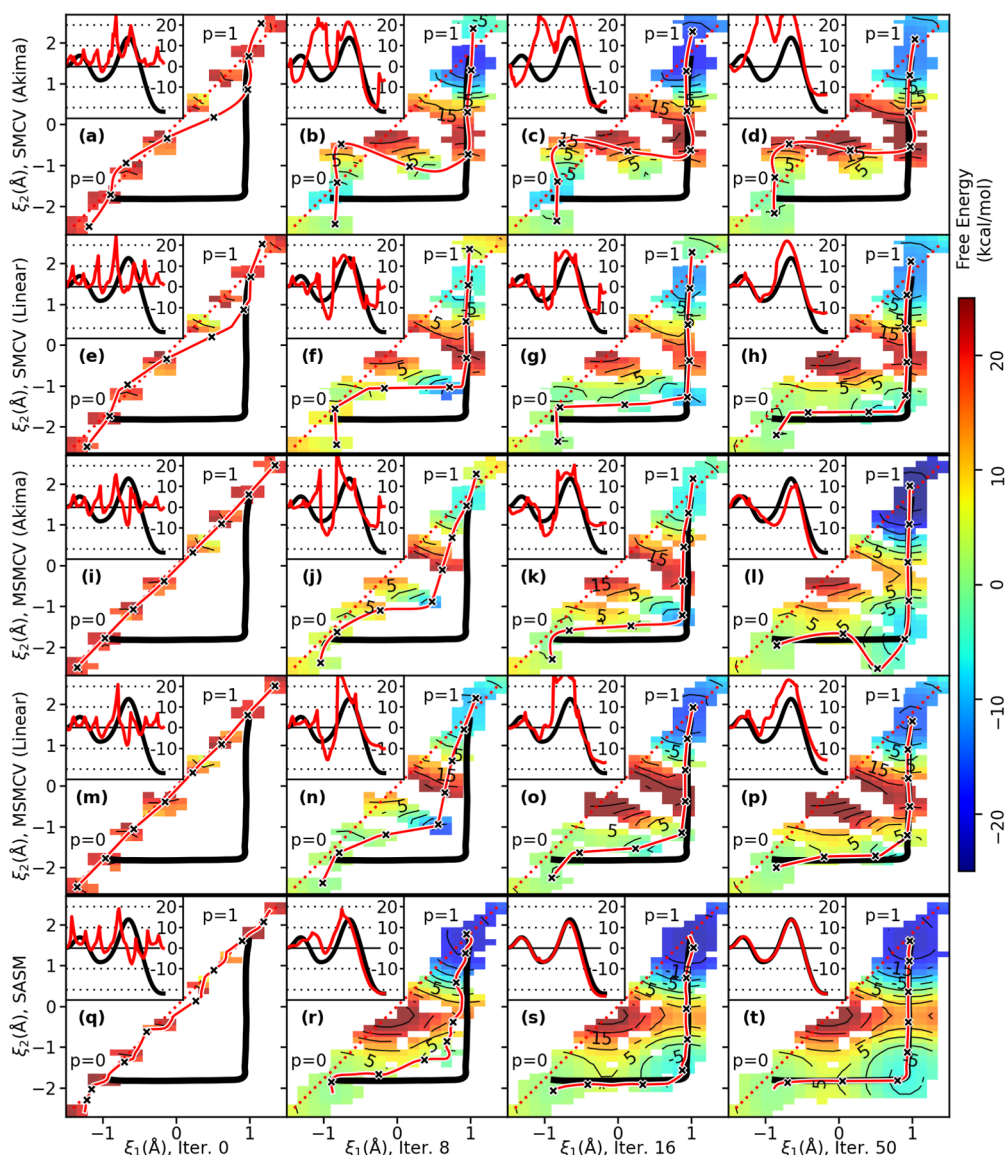


Figure 3. String iterations of MTR1 from a concerted (linear) initial guess (dashed red line). Each string is composed of 8 images, and each image is sampled for 4 ps. The solid red line is the current string, and the black “x” marks the next set of 8 simulations. The black line is a reference pathway, and the insets compare the current estimate of the free energy profile to the reference profile (kcal/mol). Parts (a–d) and (e–h) are the SMCV method using Akima and piecewise linear splines, respectively. Parts (i–l) and (m–p) similarly compare the MSMCV method. Parts (q–t) are the SASM method with 100 synthetic images.

extremely costly to perform *ab initio* QM/MM simulations. Hence, we use this as our first test system for developing improved string methods with accelerated convergence.

Figure 2 uses the MTR1 reaction to compare the progress of the SMCV, MSMCV, and SASM at string iterations 5, 15, 30, and 50. Each optimization was performed twice, starting from concerted and stepwise initial paths. Each string iteration samples 32 images for 4 ps/image (128 ps/iteration). The free energy surface is the best estimate made from the aggregate sampling of all iterations obtained from the 3 methods (38.4 ns of aggregate sampling). The black line is a reference MFEP, optimized on the aggregate free energy surface. The SMCV and MSMCV paths are Akima splines fit to the 32 evolved images, whereas the SASM paths are Akima splines fit to 100 synthetic images.

The three methods approach the MFEP at different rates. The SMCV and MSMCV make good progress during the first

15 iterations, but their progress stalls as they near the MFEP. This is due to the free energy surface becoming relatively flat near the MFEP. In contrast, the SASM gets closer to the MFEP at iteration 5 than the SMCV or MSMCV does at iteration 50. By placing the simulations around the path, the SASM is capable of exploring flat surfaces more efficiently.

To test whether the conclusions drawn from Figure 2 are sensitive to the simulation time scale (time/image), we reperformed the string methods using only 500 fs/image of sampling. The resulting comparison (Supporting Information Figure S1) is nearly indistinguishable from Figure 2.

Figure 3 compares the string methods using fewer images and different spline representations of the path. The optimizations start from a concerted path, and each iteration samples 8 images for 4 ps/image. The colored areas are the current estimate of the free energy surface from the sampling produced by the current and previous iterations. The red line is

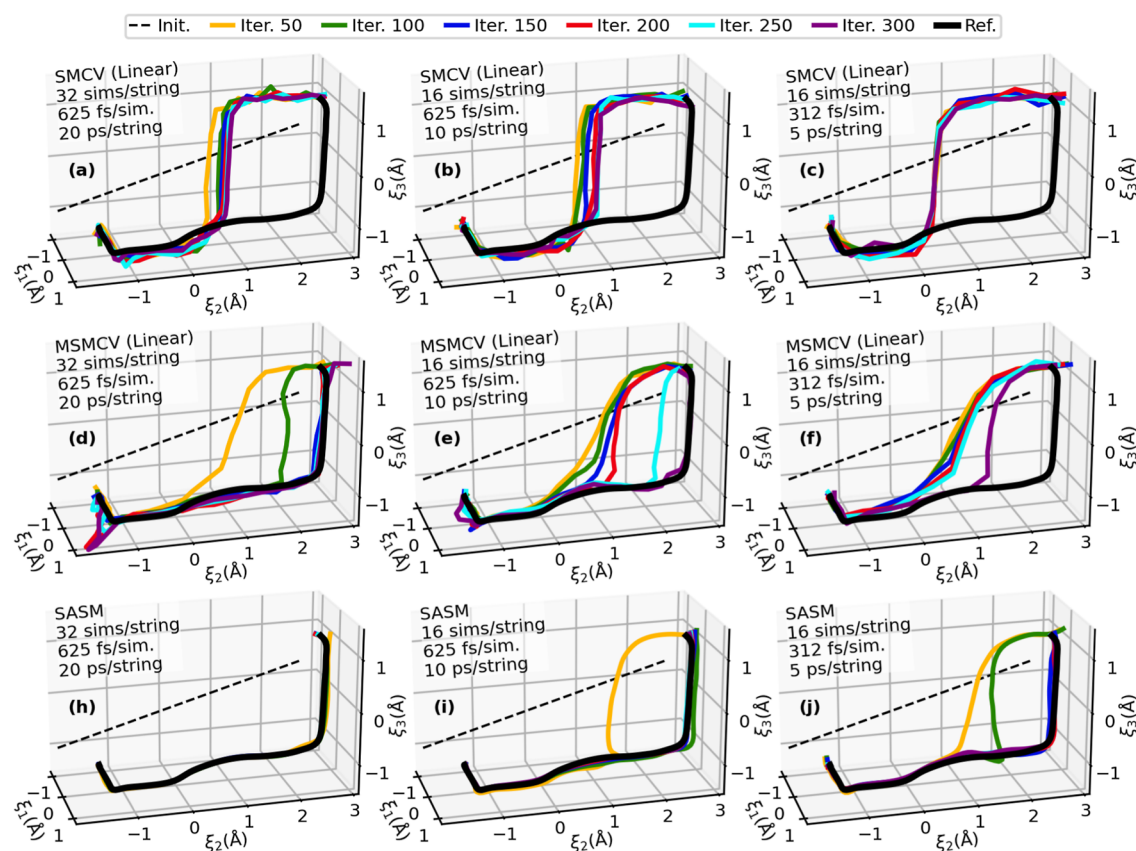


Figure 4. String iterations of HHr from a linear initial guess (dashed black line). Parts (a–c), (d–f), and (h–j) illustrate the convergence of the SMCV, MSMCV, and SASM, respectively, with several simulation protocols.

the estimate of the path after the evolution step. The “x” marks are the proposed set of umbrella potential locations. The black line is the reference path shown in Figure 2. The insets display the free energy along the paths; the red line is the free energy of the current path from the available sampling, and the black line is the reference free energy along the reference path (made from 38.4 ns of aggregate sampling). The rows labeled “Akima” and “Linear” construct the parametric curve from Akima splines and piecewise linear functions, respectively. Whereas the SMCV and MSMCV paths are limited to 8 control points, the SASM path is constructed from Akima splines which interpolate 100 synthetic images optimized on the free energy surface.

Figure 3 illustrates that the SMCV and MSMCV methods are sensitive to the parametric form of the path when only a few (e.g., 8) images are simulated. Although the SMCV and MSMCV methods properly evolve the control points, both methods encounter artifacts within the reparametrization step when the path is modeled with Akima splines. The artifacts encountered by the SMCV are quite severe; reparametrization of the curve causes some images to be propagated in a direction away from the MFEP (Figure 3a,b), and the converged path differs significantly from the reference path (Figure 3d). The MSMCV similarly encounters artifacts between iterations 15 and 50, and it converges to an incorrect path (Figure 3). The SMCV and MSMCV methods do approach the correct MFEP when using piecewise linear curves, however (see Figure 3h,p). The SASM does not exhibit artifacts using Akima splines because it is parametrized to 100 synthetic control points rather than 8 control points. Using

more control points to define the path, the SASM also avoids corner cutting, which can be observed when using piecewise linear paths; for example, see the intermediate state in Figure 3p.

When only 8 images are simulated, the progress of the SASM is modestly better than SMCV (Linear) and MSMCV (Linear); however, the SASM does a much better job at producing samples to analyze the free energy surface. As can be seen in the insets of Figures 3a–p, the limited number of images causes the SMCV and MSMCV to produce sampling that does not well overlap, resulting in noisy free energy profiles. In contrast, the SASM evolution step shifts the progress values to improve the sampling between the set of uniformly discretized points, and the exploration steps provide sampling around the path. Consequently, the SASM free energy profile after 15 iterations reproduces the reference profile very well (Figure 3s). In fact, the SASM profile after 15 iterations is better than the SMCV and MSMCV profiles after 50 iterations. The SASM placement algorithm attempts to fill the gaps in the sampling, which is easiest to observe in Figure 3q. After sampling the initial guess, the optimized path remains similar to the initial guess because all areas of the surface which have not been sampled are assumed to have a high free energy. The first 3 cases in eq 14 propose new simulations in the unoccupied regions along the path.

3.2. HHr Mechanism. The HHr^{59,102–104} is a metal-dependent small endonucleolytic self-cleaving RNA that has been extensively studied experimentally^{105–107,107,108} and computationally^{109–115} and is an archetype model for RNA catalysis. The active site adopts an L-platform/L-scaffold

architecture¹¹⁶ with an L-pocket guanine residue that forms a divalent metal ion binding site enabling electrostatic interactions⁹⁴ to facilitate the reaction. The 2'-O-transphosphorylation mechanism can be described by three reaction coordinates, illustrated in Figure 1b.

Figure 4 extends the comparisons to the 3-dimensional HHr transphosphorylation reaction profiles. The dashed line is the concerted initial guess, and the remaining lines are the paths at a series of string iterations. The reference path shown in each image is provided as a visual aid. The reference path is the SASM MFEP after 300 iterations using 32 images and 625 fs/image of sampling. In other words, it is the MFEP optimized on the 3-dimensional surface produced from the analysis of 6 ns of aggregate sampling. The string methods were performed multiple times by varying the number of images and the amount of sampling. Each column of Figure 4 successively halves the amount of sampling per string.

All of the string methods predict that the first stage of the reaction transfers a proton (the ξ_1 coordinate) from the O2' to the N1 position of the G12 general base (Figure 1b). The more interesting part of the comparison is the behavior of the paths in the ξ_2 - ξ_3 plane, where ξ_2 is the phosphoryl transfer coordinate, and ξ_3 measures the proton transfer between the O5' and the G8 general acid. The SMCV fails to locate the MFEP after 300 iterations, although it is possible that it may find the MFEP if iterated further.

The MSMCV locates the MFEP, but it requires many iterations. The MSMCV requires 150 iterations to locate the MFEP when performed with 32 images and sampled for 625 fs/image (Figure 4d). This corresponds to 3 ns of aggregate sampling. When the number of images is reduced to 16 (Figure 4e), the amount of sampling per iteration is reduced, but the MSMCV now requires 300 iterations (3 ns of aggregate sampling) to locate the MFEP. Further reduction in the amount of sampling requires more than 300 MSMCV iterations (Figure 4f). Notice that the progress of the MSMCV in Figures 4e-f does not significantly change from iterations 50 to 150, which would likely cause one to incorrectly believe that the path has converged. In fact, previous applications of the MSMCV to the HHr reaction incorrectly concluded that the mechanism was concerted because of this behavior,⁴⁹ whereas the extended iterations presented in Figure 4 suggest that the MFEP is stepwise. The fundamental reason why MSMCV progress stalls is because the free energy gradient in the directions perpendicular to the path is quite small (Figure 5). The qualitative similarity between the MSMCV path at iteration 50 to the paths produced by SMCV is suggestive that the SMCV fails for a similar reason.

In this application, the SASM requires 3 times fewer iterations than MSMCV to reach convergence when using the same amount of sampling. Only 50 SASM iterations are required to converge the path using 32 images (Figure 4h) in comparison to 150 MSMCV iterations. When the number of images is reduced to 16 (Figure 4i), convergence is reached after 100 SASM iterations in comparison to 300 MSMCV iterations. The SASM requires fewer iterations because the synthetic string optimizations performed within the SASM can evolve the path to the fringes of the aggregate sampling, and the exploration steps increase the range of the free energy surface that can be used.

3.3. B-DNA G-T Wobble Tautomer Reaction. Rare tautomeric forms of nucleobases can cause Watson-Crick-like (WC-like) mispairs in DNA, and in turn lead to disease.¹¹⁷ In

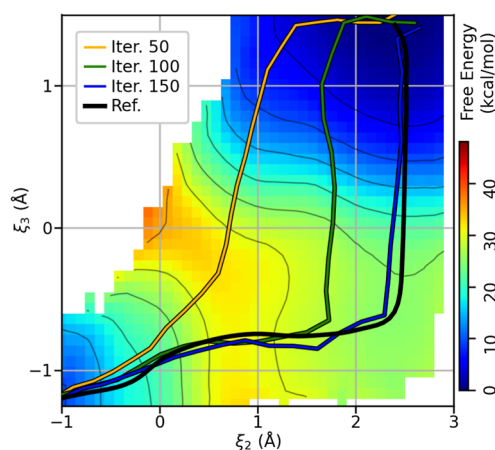


Figure 5. Two-dimensional projection of the HHr free energy surface defined by the $\xi_1 = 0.95$ Å plane. The colored lines are the MSMCV paths at iterations 50, 100, and 150, and the black line is the SASM reference curve after 300 iterations, as shown in Figure 4d. The colored areas are the free energy values calculated from the aggregate sampling produced from 300 MSMCV and 300 SASM iterations (12 ns of sampling).

the WC model, nucleobase pairs are in their “keto” form,¹¹⁸ rather than “imino” or “enol” form. Recently, tautomerization has been reported for a G-T wobble pair ($G^{\text{enol}}T/U \leftrightarrow GT^{\text{enol}}/U^{\text{enol}}$) in B-DNA detected by NMR^{97,119,120} and subsequently studied computationally.¹²¹ This tautomerization reaction can be described by 5 reaction coordinates (Figure 1c).

A pathway in 5D cannot be visualized in the same way as the 2D and 3D systems; hence, Figure 6a illustrates the convergence of the MFEP for the B-DNA tautomer reaction by calculating the root-mean-square deviation (RMSD) between the current estimate of the path and the initial guess using the 5 reaction coordinates. The SASM RMSD values plateau at 50 iterations, whereas the MSMCV requires 150 iterations to reach a similar RMSD. Although both methods seek to locate the nearest MFEP, they use different representations of the path, so one would not expect the RMSD values to exactly agree. Specifically, the MSMCV path is a piecewise linear spline constructed from 32 images, whereas the SASM path is an Akima spline constructed from 100 synthetic images. After 150 iterations, the SASM and MSMCV paths fluctuate about the MFEP; however, the fluctuations in the SASM RMSD values are significantly dampened because the additional sampling introduced by each iteration represents a smaller percentage of the aggregate. Figure 6b shows the initial and final profiles of the $\xi_1 = R_{N3-H3} - R_{O6-H3}$ and $\xi_2 = R_{O6-H3} - R_{O4-H3}$ reaction coordinates. The other 3 reaction coordinates are excluded from the figure to improve legibility. The initial path directly transfers the proton from N3 to the O4 position. The optimized paths instead transfer the proton from the N3 to O6 while shifting the hydrogen bond pattern of the G:T base pair. This is followed by the transfer of the proton from the O6 to the O4 position. The SASM and MSMCV produce very similar paths after 150 iterations. In summary, this application shows that the SASM can be extended to 5 dimensions and it can converge the path in fewer iterations than the MSMCV.

3.4. Computational Cost. Table 1 compares the CPU resources needed to perform the string methods on the HHr system with 32 images and 625 fs/image of sampling and the B-DNA system with 32 images and 1 ps/image of sampling.

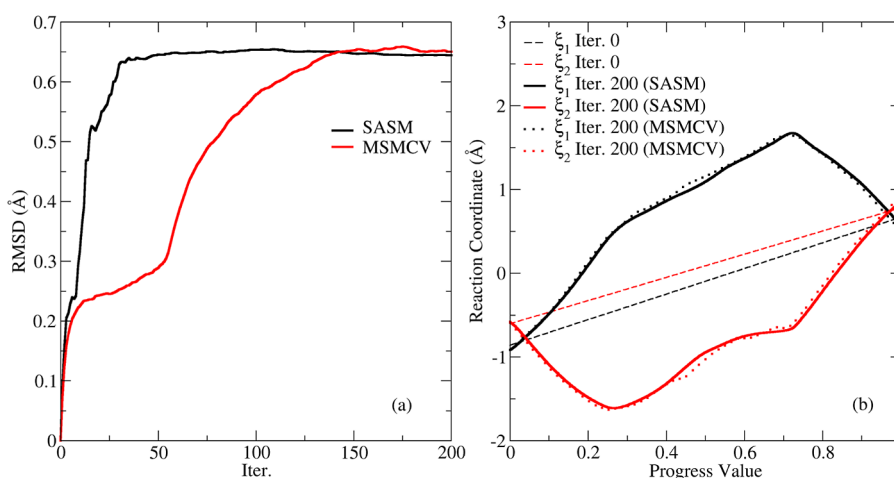


Figure 6. Convergence of the path describing the wGT → GT* tautomeric reaction in B-DNA. (a) RMSD of the 5 reaction coordinates relative to the concerted (linear) initial guess. (b) ξ_1 and ξ_2 reaction coordinates along the initial and final pathways produced by the MSMCV and SASM. These two coordinates describe the proton transfer between N3–O6 and O6–O4, respectively.

Table 1. Number of CPU Days Required to Perform MSMCV and SASM on the HHr and B-DNA Systems for the Specified Number of Iterations^a

iter	HHr			B-DNA		
	T_{MSMCV}	T_{SASM}	$T_{\text{SASM}}/T_{\text{MSMCV}}$	T_{MSMCV}	T_{SASM}	$T_{\text{SASM}}/T_{\text{MSMCV}}$
0	0.51	0.51	1.00	0.19	0.19	1.00
10	5.56	5.56	1.00	2.07	2.08	1.00
25	13.14	13.16	1.00	4.89	4.96	1.01
50	25.77	25.96	1.01	9.60	10.07	1.05
100	51.03	52.48	1.03	19.01	22.57	1.19
150	76.30	81.12	1.06	28.42	40.20	1.41

^aIteration 0 is the simulation and analysis of the initial path. Bold entries denote converged paths.

The measurements were performed on a single core of an Intel Xeon E5–2630 v3 processor, and the software was compiled with GCC 9.2.1. The timings can be decomposed into two components: the resources used to perform the QM/MM simulations T_{sim} and the resources used to perform the evolution step T_{evo} .

$$T(k) = \sum_{k'=0}^k T_{\text{sim}}(k') + T_{\text{evo}}(k') \\ = (k+1)T_{\text{sim}} + \sum_{k'=0}^k T_{\text{evo}}(k') \quad (20)$$

The MSMCV times only include the resources used to perform the QM/MM simulations; the string evolution step (eq 10) requires a negligible amount of effort, $T_{\text{evo}}(k) \approx 0$. The cost of performing MSMCV for the HHr and B-DNA systems is given by eqs 21 and 22, respectively.

$$T_{\text{MSMCV}}^{\text{HHr}}(k) = (k+1)(32 \text{ images}) \left(\frac{0.625 \text{ ps}}{\text{image}} \right) \left(\frac{1 \text{ CPU day}}{39.6 \text{ ps}} \right) \quad (21)$$

$$T_{\text{MSMCV}}^{\text{B-DNA}}(k) = (k+1)(32 \text{ images}) \left(\frac{1 \text{ ps}}{\text{image}} \right) \left(\frac{1 \text{ CPU day}}{170 \text{ ps}} \right) \quad (22)$$

The SASM timings also include the cost of the evolution step, which is further decomposed into the resources used to solve the MBAR/UWHAM equations, T_{MBAR} , and the cost of performing an optimization on the resulting free energy surface, T_{opt} .

$$T_{\text{SASM}}(k) = T_{\text{MSMCV}}(k) + \sum_{k'=0}^k T_{\text{MBAR}}(k') + T_{\text{opt}}(k') \quad (23)$$

The solution of the MBAR/UWHAM equations formally scales $O(N_{\text{dim}}N_{\text{samples}}N_{\text{states}})$, where N_{samples} is the number of samples to be reweighted and N_{states} is the number of states. The dimensionality does not vary with string iteration, and N_{samples} and N_{states} are both proportional to the number of iterations, leading to $T_{\text{MBAR}}(k) \approx A(k+1)^2$, where A is the coefficient fit to the observed times. This coefficient is 0.359 and 0.876 s for the HHr and B-DNA systems, respectively. The quadratic dependence of $T_{\text{MBAR}}(k)$ means that the aggregate cost for performing k string iterations scales cubically. The cost of performing the optimization formally scales $O(o^{N_{\text{dim}}}N_{\text{sites}}N_{\text{sim}})$, where o is the order of the Cardinal B-spline, N_{sites} is the number of synthetic iterations, and N_{sim} is the number of synthetic images used to describe the path. These quantities are independent of string iteration, so $T_{\text{opt}}(k) \approx B$, where B is 0.7 and 19.5 s for the HHr and B-DNA systems, respectively.

The timings listed in Table 1 suggest that the SASM increases the computational cost by 1–5% relative to the MSMCV for the first 50 iterations. This small increase is reflected in the high computational cost of performing QM/MM sampling. Although the SASM is more expensive, it converges in fewer iterations. The SASM reduces the resources needed to converge the path by factors of 2.9 and 2.8 for the HHr and B-DNA systems, respectively. The SASM becomes increasingly expensive with respect to the number of iterations because the MBAR/UWHAM equations are solved using aggregate sampling. To prevent the method from becoming too costly at high iterations, one could limit the analysis to the

samples produced from the most recent 50 iterations, for example.

Figure 7 uses the first 50 SASM iterations of the B-DNA system to illustrate the cost of T_{MBAR} and T_{opt} as the number of

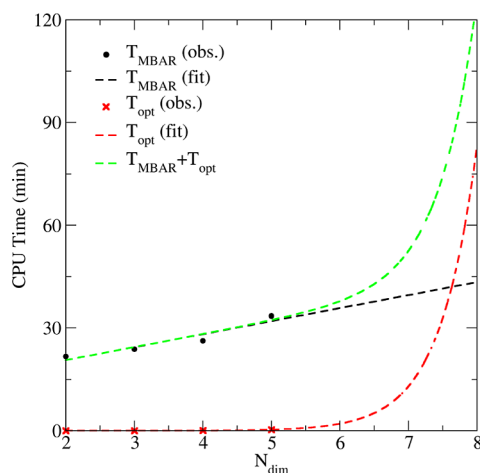


Figure 7. CPU time required to perform MBAR analysis (T_{MBAR}) and path optimization (T_{opt}) on the resulting free energy surface. The observed times were measured using the B-DNA sampling at iteration 50. The black and red dashed lines are linear and exponential fits to T_{MBAR} and T_{opt} , respectively.

reaction coordinates is varied. Although the sampling was performed with 5 reaction coordinates, we can measure T_{MBAR} and T_{opt} by ignoring 1-or-more of the reaction coordinates during the analysis. As previously discussed, the solution of the MBAR/UWHAM equations has a linear dependence on N_{dim} , and B-spline evaluations of the free energy surface have an exponential dependence on N_{dim} . The dashed lines are linear and exponential fits to the observed times. Figure 7 demonstrates that the SASM quickly becomes impractical when using more than 6 reaction coordinates due to the high cost of evaluating the free energy of a high-dimensional surface. Another aspect to consider is that each added dimension further subdivides the samples into different histogram bins. For a fixed amount of sampling, each subdivision reduces the average number of samples per occupied bin and thus increases the uncertainty of the free energy in that region. For these reasons, we do not view the SASM as a replacement for the MSMCV when a large number of reaction coordinates is needed. Instead, the SASM is a complementary tool specifically tailored to accelerate the convergence of low-dimensional pathways frequently encountered in QM/MM applications. In these situations, the added expense of generating free energy surfaces and optimizing paths from the available sampling is worthwhile to reduce the number of QM/MM evaluations.

The bond forming and breaking events of many biological mechanisms can be described with 6-or-fewer dimensions; however, this limitation may become problematic when the chemical events are coupled with conformational changes. As a specific example, a previous investigation of nucleobase tautomerization reactions used 13 interatomic distances to describe the change in hydrogen bond patterns.¹²¹ A compromise solution may be to explicitly model the bond forming and breaking coordinates with distances (or distance differences) and describe the conformational changes with path collective variables^{28,43} or other strategies developed to reduce the dimensionality.^{122–128} Alternatively, the SASM

approach could greatly benefit from new methods that are exploring the use of deep neural networks to efficiently represent high-dimensional free energy surfaces.^{129,130}

The reader may question if there are ideas introduced in the SASM that can be directly used to improve the convergence of the SMCV and MSMCV methods when a high-dimensional free energy surface is too expensive to evaluate. Unfortunately, if the free energy surface is not available, then the synthetic optimizations (eq 11) are clearly not possible. Furthermore, the SASM progress value shifting (eq 14) and exploration (eq 16) modifications rely on the fact that the simulated images are not responsible for describing the path, which is not the situation when using the SMCV or MSMCV methods. Nevertheless, the SASM approach is an enabling technology that allows one to explore new strategies that are not readily possible within the SMCV and MSMCV frameworks, as described below.

The reaction coordinates describing the chemical events are often assigned from chemical intuition, and one typically performs the string method several times starting from different initial guess pathways to compare a limited number of plausible mechanistic scenarios; for example, the associative versus dissociative mechanisms of phosphoryl transfer reactions. From a given initial guess, the SASM seeks to find the nearest MFEP. It is unlikely that the SASM will discover an alternate pathway unless it is separated by a small barrier that could be leapfrogged by the exploration stage. We have investigated the idea of using the SASM to simultaneously propagate multiple strings, each starting from a different initial guess. Each image from every string is sampled, the sampling is aggregated to form a single free energy surface, and the MFEP of each string is obtained from independent synthetic optimizations on the unified surface. In other words, the SASM is performed for each string, but their progress is synchronized at each iteration to form a single, global view of the free energy surface. The present work did not elaborate on this strategy because we remain unconvinced that it offers a meaningful computational advantage in the few test cases we have performed, which involved pathways that did not significantly overlap with each other in areas other than the reactant and product states. In these situations, it is sufficient to independently converge each string and aggregate the sampling in a postprocessing step. Furthermore, we note that the SASM formally provides the capability to sample with an inexpensive reference semiempirical Hamiltonian while propagating the string with a high-level target Hamiltonian using the weighted thermodynamic perturbation method to construct the free energy surface.^{71,131,132} We suspect, however, that a better strategy would be to converge the path with the reference potential, perform production sampling on the final path, and reweight the production sampling to estimate the target free energy surface only in the immediate vicinity of the MFEP.

4. CONCLUSIONS

We applied the SMCV, MSMCV, and SASM methods to QM/MM sampling of the MTR1, HHR, and B-DNA G-T mismatch systems. These applications served to compare the behavior and performance of the string methods using 2, 3, and 5 reaction coordinates. The SASM is a new method developed in this work that is robust and has performance advantages for systems up to approximately 6 dimensions ($N_{\text{dim}} \leq 6$). Rather than propagating the path from the sampling produced by the

most recent set of images, the SASM uses aggregate sampling from all string iterations. The sampling is used to construct the current best estimate of a multidimensional free energy surface, and an MFEP is optimized on the surface. Consequently, the simulated images are no longer responsible for describing the parametric form of the path; their sole responsibility is to improve the quality and range of the sampling used to estimate the surface. The SASM exploits this freedom by alternating between “exploration” and “refinement” steps to rapidly traverse flat regions of the free energy surface.

Overall, the SMCV, MSMCV, and SASM methods are capable of converging to the correct MFEP if the right control parameters are found. In some cases, spline artifacts can be observed with the SMCV and MSMCV when only a few images (e.g., 8) are used. The SASM is found to be more robust, and it often requires approximately 1/3 of the string iterations to converge the MFEP. Analysis of computational timings indicates that the SASM increases the computational cost per string iteration by 5% or less relative to the MSMCV, but this is more than offset by requiring fewer iterations to reach convergence. The computational cost of representing a free energy surface with more than 6 reaction coordinates quickly becomes prohibitive; therefore, the SASM is not a blanket replacement for the MSMCV. Rather, it is a valuable tool that can be used to considerably accelerate convergence in QM/MM applications using a modest number of reaction coordinates.

■ ASSOCIATED CONTENT

SI Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jctc.3c01401>.

Descriptions of the procedures used to smooth the control points and Cardinal B-spline evaluation of the free energy, and a comparison of MTR1 profiles generated from reduced sampling (PDF)

■ AUTHOR INFORMATION

Corresponding Author

Darrin M. York – Laboratory for Biomolecular Simulation Research, Institute for Quantitative Biomedicine and Department of Chemistry and Chemical Biology, Rutgers University, Piscataway, New Jersey 08854, United States; orcid.org/0000-0002-9193-7055; Email: Darrin.York@rutgers.edu

Authors

Timothy J. Giese – Laboratory for Biomolecular Simulation Research, Institute for Quantitative Biomedicine and Department of Chemistry and Chemical Biology, Rutgers University, Piscataway, New Jersey 08854, United States; orcid.org/0000-0002-0653-9168

Şölen Ekesan – Laboratory for Biomolecular Simulation Research, Institute for Quantitative Biomedicine and Department of Chemistry and Chemical Biology, Rutgers University, Piscataway, New Jersey 08854, United States; orcid.org/0000-0002-5598-5754

Erika McCarthy – Laboratory for Biomolecular Simulation Research, Institute for Quantitative Biomedicine and Department of Chemistry and Chemical Biology, Rutgers University, Piscataway, New Jersey 08854, United States; orcid.org/0000-0001-8089-0207

Yujun Tao – Laboratory for Biomolecular Simulation Research, Institute for Quantitative Biomedicine and Department of Chemistry and Chemical Biology, Rutgers University, Piscataway, New Jersey 08854, United States; orcid.org/0000-0002-4520-941X

Complete contact information is available at: <https://pubs.acs.org/10.1021/acs.jctc.3c01401>

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

The authors are grateful for the financial support provided by the National Institutes of Health (no. GM62248) and the National Science Foundation (CSSI Frameworks Grant no. 2209717). Computational resources were provided by the Office of Advanced Research Computing (OARC) at Rutgers, The State University of New Jersey; the Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support (ACCESS) program, which is supported by National Science Foundation grants #2138259, #2138286, #2138307, #2137603, and #2138296 (supercomputer Expanse at SDSC through allocation CHE190067); and the Texas Advanced Computing Center (TACC) at the University of Texas at Austin, URL: <http://www.tacc.utexas.edu> (supercomputer Frontera through allocation CHE20002).

■ REFERENCES

- (1) Giese, T. J.; York, D. M. Quantum mechanical force fields for condensed phase molecular simulations. *J. Phys.: Condens. Matter* **2017**, *29*, 383002.
- (2) Gao, J.; Truhlar, D. G. Quantum Mechanical Methods for Enzyme Kinetics. *Annu. Rev. Phys. Chem.* **2002**, *53*, 467–505.
- (3) Garcia-Viloca, M.; Gao, J.; Karplus, M.; Truhlar, D. G. How enzymes work: Analysis by modern rate theory and computer simulations. *Science* **2004**, *303*, 186–195.
- (4) Giese, T. J.; Huang, M.; Chen, H.; York, D. M. Recent Advances toward a General Purpose Linear-Scaling Quantum Force Field. *Acc. Chem. Res.* **2014**, *47*, 2812–2820.
- (5) Gao, J.; Truhlar, D. G.; Wang, Y.; Mazack, M. J.; Löffler, P.; Provorse, M. R.; Rehak, P. Explicit polarization: A quantum mechanical framework for developing next generation force fields. *Acc. Chem. Res.* **2014**, *47*, 2837–2845.
- (6) Behler, J. First Principles Neural Network Potentials for Reactive Simulations of Large Molecular and Condensed Systems. *Angew. Chem., Int. Ed.* **2017**, *56*, 12828–12840.
- (7) Meuwly, M. Machine Learning for Chemical Reactions. *Chem. Rev.* **2021**, *121*, 10218–10239.
- (8) Zeng, J.; Giese, T. J.; Ekesan, S.; York, D. M. Development of Range-Corrected Deep Learning Potentials for Fast, Accurate Quantum Mechanical/Molecular Mechanical Simulations of Chemical Reactions in Solution. *J. Chem. Theory Comput.* **2021**, *17*, 6993–7009.
- (9) Giese, T. J.; Zeng, J.; Ekesan, S.; York, D. M. Combined QM/MM, Machine Learning Path Integral Approach to Compute Free Energy Profiles and Kinetic Isotope Effects in RNA Cleavage Reactions. *J. Chem. Theory Comput.* **2022**, *18*, 4304–4317.
- (10) Pan, X.; Yang, J.; Van, R.; Epifanovsky, E.; Ho, J.; Huang, J.; Pu, J.; Mei, Y.; Nam, K.; Shao, Y. Machine-Learning-Assisted Free Energy Simulation of Solution-Phase and Enzyme Reactions. *J. Chem. Theory Comput.* **2021**, *17*, 5745–5758.
- (11) Snyder, R.; Kim, B.; Pan, X.; Shao, Y.; Pu, J. Bridging semiempirical and *ab initio* QM/MM potentials by Gaussian process regression and its sparse variants for free energy simulation. *J. Chem. Phys.* **2023**, *159*, 054107.

- (12) Kästner, J. Umbrella sampling. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **2011**, *1*, 932–942.
- (13) Torrie, G. M.; Valleau, J. P. Monte Carlo free energy estimates using non-Boltzmann sampling: Application to the sub-critical Lennard-Jones fluid. *Chem. Phys. Lett.* **1974**, *28*, 578–581.
- (14) Torrie, G. M.; Valleau, J. P. Nonphysical sampling distributions in Monte Carlo free-energy estimation: Umbrella sampling. *J. Comput. Phys.* **1977**, *23*, 187–199.
- (15) McDonald, I. R.; Singer, K. Machine Calculation of Thermodynamic Properties of a Simple Fluid at Supercritical Temperatures. *J. Chem. Phys.* **1967**, *47*, 4766–4772.
- (16) McDonald, I. R.; Singer, K. Examination of the Adequacy of the 12–6 Potential for Liquid Argon by Means of Monte Carlo Calculations. *J. Chem. Phys.* **1969**, *50*, 2308–2315.
- (17) Adib, A. B. Free energy surfaces from nonequilibrium processes without work measurement. *J. Chem. Phys.* **2006**, *124*, 144111.
- (18) Hummer, G. Fast-growth thermodynamic integration: Error and efficiency analysis. *J. Chem. Phys.* **2001**, *114*, 7330–7337.
- (19) Zuckerman, D. M.; Woolf, T. B. Theory of a Systematic Computational Error in Free Energy Differences. *Phys. Rev. Lett.* **2002**, *89*, 180602.
- (20) Ozer, G.; Valeev, E. F.; Quirk, S.; Hernandez, R. Adaptive Steered Molecular Dynamics of the Long-Distance Unfolding of Neuropeptide Y. *J. Chem. Theory Comput.* **2010**, *6*, 3026–3038.
- (21) Jarzynski, C. Nonequilibrium equality for free energy differences. *Phys. Rev. Lett.* **1997**, *78*, 2690–2693.
- (22) Knight, J. L.; Brooks, C. L. λ -Dynamics free energy simulation methods. *J. Comput. Chem.* **2009**, *30*, 1692–1700.
- (23) Kong, X.; Brooks, C. L. λ -dynamics: A new approach to free energy calculations. *J. Chem. Phys.* **1996**, *105*, 2414–2423.
- (24) Liu, Z.; Berne, B. J. Method for accelerating chain folding and mixing. *J. Chem. Phys.* **1993**, *99*, 6071–6077.
- (25) Tidor, B. Simulated annealing on free energy surfaces by a combined molecular dynamics and Monte Carlo approach. *J. Phys. Chem.* **1993**, *97*, 1069–1073.
- (26) Laio, A.; Parrinello, M. Escaping free-energy minima. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, *99*, 12562–12566.
- (27) White, A. D.; Dama, J. F.; Voth, G. A. Designing Free Energy Surfaces That Match Experimental Data with Metadynamics. *J. Chem. Theory Comput.* **2015**, *11*, 2451–2460.
- (28) Branduardi, D.; Gervasio, F. L.; Parrinello, M. From A to B in free energy space. *J. Chem. Phys.* **2007**, *126*, 054103–054112.
- (29) Sugita, Y.; Okamoto, Y. Replica-exchange molecular dynamics method for protein folding. *Chem. Phys. Lett.* **1999**, *314*, 141–151.
- (30) Yang, L.; Shao, Q.; Gao, Y. Comparison between integrated and parallel tempering methods in enhanced sampling simulations. *J. Chem. Phys.* **2009**, *130*, 124111.
- (31) Jónsson, H.; Mills, G.; Jacobsen, K. W. *Classical and Quantum Dynamics in Condensed Phase Simulations*; World Scientific, 1998; pp 385–404, Chapter Nudged elastic band method for finding minimum energy paths of transitions.
- (32) E, W.; Ren, W.; Vanden-Eijnden, E. String method for the study of rare events. *Phys. Rev. B: Condens. Matter Mater. Phys.* **2002**, *66*, 052301.
- (33) E, W.; Ren, W.; Vanden-Eijnden, E. Simplified and improved string method for computing the minimum energy paths in barrier-crossing events. *J. Chem. Phys.* **2007**, *126*, 164103.
- (34) Maragliano, L.; Fischer, A.; Vanden-Eijnden, E.; Ciccotti, G. String method in collective variables: minimum free energy paths and isocommittor surfaces. *J. Chem. Phys.* **2006**, *125*, 024106.
- (35) Ovchinnikov, V.; Karplus, M.; Vanden-Eijnden, E. Free energy of conformational transition paths in biomolecules: The string method and its application to myosin VI. *J. Chem. Phys.* **2011**, *134*, 085103.
- (36) Ovchinnikov, V.; Karplus, M. Investigations of α -helix \leftrightarrow β -sheet transition pathways in a miniprotein using the finite-temperature string method. *J. Chem. Phys.* **2014**, *140*, 175103–175121.
- (37) Vanden-Eijnden, E.; Venturoli, M. Revisiting the finite temperature string method for the calculation of reaction tubes and free energies. *J. Chem. Phys.* **2009**, *130*, 194103.
- (38) Rosta, E.; Nowotny, M.; Yang, W.; Hummer, G. Catalytic Mechanism of RNA Backbone Cleavage by Ribonuclease H from Quantum Mechanics/Molecular mechanics simulations. *J. Am. Chem. Soc.* **2011**, *133*, 8934–8941.
- (39) Khavrutskii, I. V.; Arora, K.; Brooks, C. L. Harmonic Fourier beads method for studying rare events on rugged energy surfaces. *J. Chem. Phys.* **2006**, *125*, 174108.
- (40) Peters, B.; Heyden, A.; Bell, A. T.; Chakraborty, A. A growing string method for determining transition states: comparison to the nudged elastic band and string methods. *J. Chem. Phys.* **2004**, *120*, 7877–7886.
- (41) Pan, A. C.; Sezer, D.; Roux, B. Finding transition pathways using the string method with swarms of trajectories. *J. Phys. Chem. B* **2008**, *112*, 3432–3440.
- (42) Matsunaga, Y.; Komuro, Y.; Kobayashi, C.; Jung, J.; Mori, T.; Sugita, Y. Dimensionality of Collective Variables for Describing Conformational Changes of a Multi-Domain Protein. *J. Phys. Chem. Lett.* **2016**, *7*, 1446–1451.
- (43) Kulshrestha, A.; Punnathanam, S. N.; Ayappa, K. G. Finite temperature string method with umbrella sampling using path collective variables: application to secondary structure change in a protein. *Soft Matter* **2022**, *18*, 7593–7603.
- (44) Ensing, B.; Laio, A.; Parrinello, M.; Klein, M. L. A Recipe for the Computation of the Free Energy Barrier and the Lowest Free Energy Path of Concerted Reactions. *J. Phys. Chem. B* **2005**, *109*, 6676–6687.
- (45) Badaoui, M.; Buigues, P. J.; Berta, D.; Mandana, G. M.; Gu, H.; Földes, T.; Dickson, C. J.; Hornak, V.; Kato, M.; Molteni, C.; Parsons, S.; Rosta, E. Combined Free-Energy Calculation and Machine Learning Methods for Understanding Ligand Unbinding Kinetics. *J. Chem. Theory Comput.* **2022**, *18*, 2543–2555.
- (46) Ganguly, A.; Thaplyal, P.; Rosta, E.; Bevilacqua, P. C.; Hammes-Schiffer, S. Quantum Mechanical/Molecular Mechanical Free Energy Simulations of the Self-Cleavage Reaction in the Hepatitis Delta Virus Ribozyme. *J. Am. Chem. Soc.* **2014**, *136*, 1483–1496.
- (47) Zhong, J.; Reinhardt, C. R.; Hammes-Schiffer, S. Role of Water in Proton-Coupled Electron Transfer between Tyrosine and Cysteine in Ribonucleotide Reductase. *J. Am. Chem. Soc.* **2022**, *144*, 7208–7214.
- (48) Zhong, J.; Reinhardt, C. R.; Hammes-Schiffer, S. Direct Proton-Coupled Electron Transfer between Interfacial Tyrosines in Ribonucleotide Reductase. *J. Am. Chem. Soc.* **2023**, *145*, 4784–4790.
- (49) Giese, T. J.; Ekesan, S.; York, D. M. Extension of the Variational Free Energy Profile and Multistate Bennett Acceptance Ratio Methods for High-Dimensional Potential of Mean Force Profile Analysis. *J. Phys. Chem. A* **2021**, *125*, 4216–4232.
- (50) Reinhardt, C. R.; Sayfutyarova, E. R.; Zhong, J.; Hammes-Schiffer, S. Glutamate Mediates Proton-Coupled Electron Transfer Between Tyrosines 730 and 731 in *Escherichia coli* Ribonucleotide Reductase. *J. Am. Chem. Soc.* **2021**, *143*, 6054–6059.
- (51) Li, P.; Soudackov, A. V.; Hammes-Schiffer, S. Fundamental Insights into Proton-Coupled Electron Transfer in Soybean Lipoygenase from Quantum Mechanical/Molecular Mechanical Free Energy Simulations. *J. Am. Chem. Soc.* **2018**, *140*, 3068–3076.
- (52) Ganguly, A.; Boulanger, E.; Thiel, W. Importance of MM Polarization in QM/MM Studies of Enzymatic Reactions: Assessment of the QM/MM Drude Oscillator Model. *J. Chem. Theory Comput.* **2017**, *13*, 2954–2961.
- (53) Giese, T. J.; York, D. M. *FE-Toolkit: The Free Energy Analysis Toolkit*; Laboratory for Biomolecular Simulation Research. <https://gitlab.com/RutgersLBSR/fe-toolkit>, 2023.
- (54) Case, D. A.; Aktulga, H. M.; Belfon, K.; Cerutti, D. S.; Cisneros, G. A.; Cruzeiro, V. W. D.; Forouzeshe, N.; Giese, T. J.; Götz, A. W.; Gohlke, H.; Izadi, S.; Kasavajhala, K.; Kaymak, M. C.; King, E.; Kurtzman, T.; Lee, T.-S.; Li, P.; Liu, J.; Luchko, T.; Luo, R.;

- Manathunga, M.; Machado, M. R.; Nguyen, H. M.; O'Hearn, K. A.; Onufriev, A. V.; Pan, F.; Pantano, S.; Qi, R.; Rahnamoun, A.; Risheh, A.; Schott-Verdugo, S.; Shajan, A.; Swails, J.; Wang, J.; Wei, H.; Wu, X.; Wu, Y.; Zhang, S.; Zhao, S.; Zhu, Q.; Cheatham, T. E.; Roe, D. R.; Roitberg, A.; Simmerling, C.; York, D. M.; Nagan, M. C.; Merz, K. M. AmberTools. *J. Chem. Inf. Model.* **2023**, *63*, 6183–6191.
- (55) Wojtas-Niziurski, W.; Meng, Y.; Roux, B.; Bernèche, S. Self-learning adaptive umbrella sampling method for the determination of free energy landscapes in multiple dimensions. *J. Chem. Theory Comput.* **2013**, *9*, 1885–1895.
- (56) Scheitl, C. P. M.; Ghaem Maghami, M.; Lenz, A.-K.; Höbartner, C. Site-specific RNA methylation by a methyltransferase ribozyme. *Nature* **2020**, *587*, 663–667.
- (57) Scheitl, C. P. M.; Mieczkowski, M.; Schindelin, H.; Höbartner, C. Structure and mechanism of the methyltransferase ribozyme MTR1. *Nat. Chem. Biol.* **2022**, *18*, 547–555.
- (58) Deng, J.; Wilson, T. J.; Wang, J.; Peng, X.; Li, M.; Lin, X.; Liao, W.; Lilley, D. M. J.; Huang, L. Structure and mechanism of a methyltransferase ribozyme. *Nat. Chem. Biol.* **2022**, *18*, 556–564.
- (59) Martick, M.; Lee, T.-S.; York, D. M.; Scott, W. G. Solvent structure and hammerhead ribozyme catalysis. *Chem. Biol.* **2008**, *15*, 332–342.
- (60) Hunter, W. N.; Brown, T.; Kneale, G.; Anand, N. N.; Rabinovich, D.; Kennard, O. The structure of guanosine-thymidine mismatches in B-DNA at 2.5-Å resolution. *J. Biol. Chem.* **1987**, *262*, 9962–9970.
- (61) Wilson, E. B.; Decius, J. C.; Cross, P. C. *Molecular Vibrations*; Dover Publications, Inc.: New York, 1980.
- (62) Akima, H. A. A New Method of Interpolation and Smooth Curve Fitting Based on Local Procedures. *J. ACM* **1970**, *17*, 589–602.
- (63) Lee, T.-S.; Radak, B. K.; Pabis, A.; York, D. M. A new maximum likelihood approach for free energy profile construction from molecular simulations. *J. Chem. Theory Comput.* **2013**, *9*, 153–164.
- (64) Lee, T.-S.; Radak, B. K.; Huang, M.; Wong, K.-Y.; York, D. M. Roadmaps through free energy landscapes calculated using the multidimensional vFEP approach. *J. Chem. Theory Comput.* **2014**, *10*, 24–34.
- (65) Shirts, M. R.; Chodera, J. D. Statistically optimal analysis of samples from multiple equilibrium states. *J. Chem. Phys.* **2008**, *129*, 124105.
- (66) Tan, Z.; Gallicchio, E.; Lapelosa, M.; Levy, R. M. Theory of binless multi-state free energy estimation with applications to protein-ligand binding. *J. Chem. Phys.* **2012**, *136*, 144102.
- (67) Zhang, B. W.; Xia, J.; Tan, Z.; Levy, R. M. A Stochastic Solution to the Unbinned WHAM Equations. *J. Phys. Chem. Lett.* **2015**, *6*, 3834–3840.
- (68) Milovanović, G. V.; Udovičić, Z. Calculation of coefficients of a cardinal B-spline. *Appl. Math. Lett.* **2010**, *23*, 1346–1350.
- (69) Hardy, R. L. Multiquadric equations of topography and other irregular surfaces. *J. Geophys. Res.* **1971**, *76*, 1905–1915.
- (70) Fornberg, B.; Wright, G. Stable computation of multiquadric interpolants for all values of the shape parameter. *Comput. Math. Appl.* **2004**, *48*, 853–867.
- (71) Li, P.; Jia, X.; Pan, X.; Shao, Y.; Mei, Y. Accelerated Computation of Free Energy Profile at ab Initio Quantum Mechanical/Molecular Mechanics Accuracy via a Semi-Empirical Reference Potential. I. Weighted Thermodynamics Perturbation. *J. Chem. Theory Comput.* **2018**, *14*, 5583–5596.
- (72) Zhang, Z.; Liu, X.; Chen, Z.; Zheng, H.; Yan, K.; Liu, J. A unified thermostat scheme for efficient configurational sampling for classical/quantum canonical ensembles via molecular dynamics. *J. Chem. Phys.* **2017**, *147*, 034109.
- (73) Zhang, Z.; Liu, X.; Yan, K.; Tuckerman, M. E.; Liu, J. Unified Efficient Thermostat Scheme for the Canonical Ensemble with Holonomic or Isokinetic Constraints via Molecular Dynamics. *J. Phys. Chem. A* **2019**, *123*, 6056–6079.
- (74) Ryckaert, J. P.; Ciccotti, G.; Berendsen, H. J. C. Numerical Integration of the Cartesian Equations of Motion of a System with Constraints: Molecular Dynamics of n-Alkanes. *J. Comput. Phys.* **1977**, *23*, 327–341.
- (75) Warshel, A.; Levitt, M. Theoretical studies of enzymic reactions: Dielectric, electrostatic and steric stabilization of the carbonium ion in the reaction of lysozyme. *J. Mol. Biol.* **1976**, *103*, 227–249.
- (76) Singh, U. C.; Kollman, P. A. A combined *ab initio* quantum mechanical and molecular mechanical method for carrying out simulations on complex molecular systems: Applications to the CH₃Cl+Cl⁻ exchange reaction and gas phase protonation of polyethers. *J. Comput. Chem.* **1986**, *7*, 718–730.
- (77) Essmann, U.; Perera, L.; Berkowitz, M. L.; Darden, T.; Lee, H.; Pedersen, L. G. A smooth particle mesh Ewald method. *J. Chem. Phys.* **1995**, *103*, 8577–8593.
- (78) Darden, T.; York, D.; Pedersen, L. Particle mesh Ewald: An N log(N) method for Ewald sums in large systems. *J. Chem. Phys.* **1993**, *98*, 10089–10092.
- (79) Petersen, H. G. Accuracy and efficiency of the particle mesh Ewald method. *J. Chem. Phys.* **1995**, *103*, 3668–3679.
- (80) Nam, K.; Gao, J.; York, D. M. An efficient linear-scaling Ewald method for long-range electrostatic interactions in combined QM/MM calculations. *J. Chem. Theory Comput.* **2005**, *1*, 2–13.
- (81) Walker, R. C.; Crowley, M. F.; Case, D. A. The implementation of a fast and accurate QM/MM potential method in Amber. *J. Comput. Chem.* **2008**, *29*, 1019–1031.
- (82) Figueirido, F.; Del Buono, G. S.; Levy, R. M. On finite-size effects in computer simulations using the Ewald potential. *J. Chem. Phys.* **1995**, *103*, 6133–6142.
- (83) de Leeuw, S. W.; Perram, J. W.; Smith, E. R. Simulation of electrostatic systems in periodic boundary conditions. I. Lattice sums and dielectric constants. *Proc. R. Soc. London, Ser. A* **1980**, *373*, 27–56.
- (84) Giese, T. J.; York, D. M. A GPU-Accelerated Parameter Interpolation Thermodynamic Integration Free Energy Method. *J. Chem. Theory Comput.* **2018**, *14*, 1564–1582.
- (85) Zgarbová, M.; Otyepka, M.; Šponer, J.; Mládek, A.; Banáš, P.; Cheatham, T. E.; Jurečka, P. Refinement of the Cornell et al. nucleic acids force field based on reference quantum chemical calculations of glycosidic torsion profiles. *J. Chem. Theory Comput.* **2011**, *7*, 2886–2902.
- (86) Joung, I. S.; Cheatham, T. E. Determination of alkali and halide monovalent ion parameters for use in explicitly solvated biomolecular simulations. *J. Phys. Chem. B* **2008**, *112*, 9020–9041.
- (87) McCarthy, E.; Ekesan, S.; Giese, T. J.; Wilson, T. J.; Deng, J.; Huang, L.; Lilley, D. J.; York, D. M. Catalytic mechanism and pH dependence of a methyltransferase ribozyme (MTR1) from computational enzymology. *Nucleic Acids Res.* **2023**, *51*, 4508–4518.
- (88) Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.; Dinola, A.; Haak, J. R. Molecular dynamics with coupling to an external bath. *J. Chem. Phys.* **1984**, *81*, 3684–3690.
- (89) Loncharich, R. J.; Brooks, B. R.; Pastor, R. W. Langevin dynamics of peptides: the frictional dependence of isomerization rates of N-acetylalanine-N'-methylamide. *Biopolymers* **1992**, *32*, 523–535.
- (90) Gaus, M.; Cui, Q.; Elstner, M. DFTB3: Extension of the Self-Consistent-Charge Density-Functional Tight-Binding Method (SCC-DFTB). *J. Chem. Theory Comput.* **2011**, *7*, 931–948.
- (91) Li, P.; Roberts, B. P.; Chakravorty, D. K.; Merz, K. M. Rational design of Particle Mesh Ewald compatible Lennard-Jones parameters for + 2 metal cations in explicit solvent. *J. Chem. Theory Comput.* **2013**, *9*, 2733–2748.
- (92) Panteva, M. T.; Giambaşu, G. M.; York, D. M. Comparison of structural, thermodynamic, kinetic and mass transport properties of Mg²⁺ ion models commonly used in biomolecular simulations. *J. Comput. Chem.* **2015**, *36*, 970–982.
- (93) Panteva, M. T.; Giambaşu, G. M.; York, D. M. Force Field for Mg²⁺, Mn²⁺, Zn²⁺, and Cd²⁺ Ions that have Balanced Interactions with Nucleic Acids. *J. Phys. Chem. B* **2015**, *119*, 15460–15470.

- (94) Ekesan, S. J.; McCarthy, E.; Case, D. A.; York, D. M. RNA Electrostatics: How Ribozymes Engineer Active Sites to Enable Catalysis. *J. Phys. Chem. B* **2022**, *126*, 5982–5990.
- (95) Lopez, X.; York, D. M. Parameterization of semiempirical methods to treat nucleophilic attacks to biological phosphates: AM1/d parameters for phosphorus. *Theor. Chem. Acc.* **2003**, *109*, 149–159.
- (96) Zgarbova, M.; Sponer, J.; Otyepka, M.; Cheatham, T. E.; Galindo-Murillo, R.; Jurecka, P. Refinement of the sugar-phosphate backbone torsion beta for AMBER Force Fields improves the description of Z- and B-DNA. *J. Chem. Theory Comput.* **2015**, *11*, 5723–5736.
- (97) Kimsey, I. J.; Szymanski, E. S.; Zahurancik, W. J.; Shakya, A.; Xue, Y.; Chu, C.-C.; Sathyamoorthy, B. J.; Suo, Z.; Al-Hashimi, H. M. Dynamic basis for dG:dT misincorporation via tautomerization and ionization. *Nature* **2018**, *554*, 195–201.
- (98) Kruger, K.; Grabowski, P.; Zaug, A.; Sands, J.; Gottschling, D.; Cech, T. Self-splicing RNA: autoexcision and autocyclisation of the ribosomal RNA intervening sequence of *Tetrahymena*. *Cell* **1982**, *31*, 147–157.
- (99) Gilbert, W. Origin of life: The RNA world. *Nature* **1986**, *319*, 618.
- (100) Deng, J.; Shi, Y.; Peng, X.; He, Y.; Chen, X.; Li, M.; Lin, X.; Liao, W.; Huang, Y.; Jiang, T.; Lilley, D. J.; Miao, Z.; Huang, L. Ribocentre: a database of ribozymes. *Nucleic Acids Res.* **2023**, *51*, D262–D268.
- (101) Wilson, T. J.; Lilley, D. M. J. The potential versatility of RNA catalysis. *Wiley Interdiscip. Rev.: RNA* **2021**, *12*, 1651.
- (102) Scott, W. G. Biophysical and biochemical investigations of RNA catalysis in the hammerhead ribozyme. *Q. Rev. Biophys.* **1999**, *32*, 241–284.
- (103) Martick, M.; Scott, W. G. Tertiary contacts distant from the active site prime a ribozyme for catalysis. *Cell* **2006**, *126*, 309–320.
- (104) Leclerc, F. Hammerhead Ribozymes: True Metal or Nucleobase Catalysis? Where Is the Catalytic Power from? *Molecules* **2010**, *15*, 5389–5407.
- (105) Blount, K. F.; Uhlenbeck, O. C. The Structure-Function Dilemma of the Hammerhead Ribozyme. *Annu. Rev. Biophys. Biomol. Struct.* **2005**, *34*, 415–440.
- (106) Nelson, J. A.; Uhlenbeck, O. C. Hammerhead redux: does the new structure fit the old biochemical data? *RNA* **2008**, *14*, 605–615.
- (107) Thomas, J. M.; Perrin, D. M. Probing general acid catalysis in the hammerhead ribozyme. *J. Am. Chem. Soc.* **2009**, *131*, 1135–1143.
- (108) Ward, W. L.; Plakos, K.; DeRose, V. J. Nucleic acid catalysis: metals, nucleobases, and other cofactors. *Chem. Rev.* **2014**, *114*, 4318–4342.
- (109) Lee, T.-S.; Silva-Lopez, C.; Martick, M.; Scott, W. G.; York, D. M. Insight into the role of Mg²⁺ in hammerhead ribozyme catalysis from x-ray crystallography and molecular dynamics simulation. *J. Chem. Theory Comput.* **2007**, *3*, 325–327.
- (110) Lee, T.-S.; López, C. S.; Giambaşu, G. M.; Martick, M.; Scott, W. G.; York, D. M. Role of Mg²⁺ in hammerhead ribozyme catalysis from molecular simulation. *J. Am. Chem. Soc.* **2008**, *130*, 3053–3064.
- (111) Lee, T.-S.; Giambaşu, G. M.; Sosa, C. P.; Martick, M.; Scott, W. G.; York, D. M. Threshold Occupancy and Specific Cation Binding Modes in the Hammerhead Ribozyme Active Site are Required for Active Conformation. *J. Mol. Biol.* **2009**, *388*, 195–206.
- (112) Lee, T.-S.; York, D. M. Computational mutagenesis studies of hammerhead ribozyme catalysis. *J. Am. Chem. Soc.* **2010**, *132*, 13505–13518.
- (113) Wong, K.-Y.; Lee, T.-S.; York, D. M. Active participation of the Mg²⁺ ion in the reaction coordinate of RNA self-cleavage catalyzed by the hammerhead ribozyme. *J. Chem. Theory Comput.* **2011**, *7*, 1–3.
- (114) Chen, H.; Giese, T. J.; Golden, B. L.; York, D. M. Divalent Metal Ion Activation of a Guanine General Base in the Hammerhead Ribozyme: Insights from Molecular Simulations. *Biochemistry* **2017**, *56*, 2985–2994.
- (115) Ganguly, A.; Weissman, B. P.; Piccirilli, J. A.; York, D. M. Evidence for a Catalytic Strategy to Promote Nucleophile Activation in Metal-Dependent RNA-Cleaving Ribozymes and 8–17 DNAzyme. *ACS Catal.* **2019**, *9*, 10612–10617.
- (116) Gaines, C. S.; Piccirilli, J. A.; York, D. M. The L-platform/L-scaffold framework: a blueprint for RNA-cleaving nucleic acid enzyme design. *RNA* **2020**, *26*, 111–125.
- (117) Mendiratta, G.; Ke, E.; Aziz, M.; Liarakos, D.; Tong, M.; SITES, E. C. Cancer gene mutation frequencies for the U.S. population. *Nat. Commun.* **2021**, *12*, 5961.
- (118) Watson, J. D.; Crick, F. H. C. Molecular Structure of Nucleic Acids: A Structure for Deoxyribose Nucleic Acid. *Nature* **1953**, *171*, 737–738.
- (119) Kimsey, I. J.; Petzold, K.; Sathyamoorthy, B.; Stein, Z. W.; Al-Hashimi, H. M. Visualizing transient Watson–Crick-like mispairs in DNA and RNA duplexes. *Nature* **2015**, *519*, 315–320.
- (120) Szymanski, E. S.; Kimsey, I. J.; Al-Hashimi, H. M. Direct NMR Evidence that Transient Tautomeric and Anionic States in dG:dT Form Watson–Crick-like Base Pairs. *J. Am. Chem. Soc.* **2017**, *139*, 4326–4329.
- (121) Li, P.; Rangadurai, A.; Al-Hashimi, H. M.; Hammes-Schiffer, S. Environmental Effects on Guanine–Thymine Mispair Tautomerization Explored with Quantum Mechanical/Molecular Mechanical Free Energy Simulations. *J. Am. Chem. Soc.* **2020**, *142*, 11183–11191.
- (122) Ceriotti, M.; Tribello, G. A.; Parrinello, M. Simplifying the representation of complex free-energy landscapes using sketch-map. *Proc. Natl. Acad. Sci. U.S.A.* **2011**, *108*, 13023–13028.
- (123) Ferguson, A. L.; Panagiotopoulos, A. Z.; Debenedetti, P. G.; Kevrekidis, I. G. Integrating diffusion maps with umbrella sampling: Application to alanine dipeptide. *J. Chem. Phys.* **2011**, *134*, 135103.
- (124) Noé, F.; Clementi, C. Kinetic Distance and Kinetic Maps from Molecular Dynamics Simulation. *J. Chem. Theory Comput.* **2015**, *11*, 5002–5011.
- (125) Tiwary, P.; Berne, B. J. Spectral gap optimization of order parameters for sampling complex molecular systems. *Proc. Natl. Acad. Sci. U.S.A.* **2016**, *113*, 2839–2844.
- (126) Das, P.; Moll, M.; Stamati, H.; Kaviraki, L. E.; Clementi, C. Low-dimensional, free-energy landscapes of protein-folding reactions by nonlinear dimensionality reduction. *Proc. Natl. Acad. Sci. U.S.A.* **2006**, *103*, 9885–9890.
- (127) Mendels, D.; Piccini, G.; Parrinello, M. Collective Variables from Local Fluctuations. *J. Phys. Chem. Lett.* **2018**, *9*, 2776–2781.
- (128) Rohrdanz, M. A.; Zheng, W.; Maggioni, M.; Clementi, C. Determination of reaction coordinates via locally scaled diffusion map. *J. Chem. Phys.* **2011**, *134*, 124116.
- (129) Zhang, L.; Wang, H.; E, W. Reinforced dynamics for enhanced sampling in large atomic and molecular systems. *J. Chem. Phys.* **2018**, *148*, 124113.
- (130) Schneider, E.; Dai, L.; Topper, R. Q.; Drechsel-Grau, C.; Tuckerman, M. E. Stochastic Neural Network Approach for Learning High-Dimensional Free Energy Surfaces. *Phys. Rev. Lett.* **2017**, *119*, 150601.
- (131) Giese, T. J.; Zeng, J.; York, D. M. Multireference Generalization of the Weighted Thermodynamic Perturbation Method. *J. Phys. Chem. A* **2022**, *126*, 8519–8533.
- (132) Wang, J.-N.; Xue, Y.; Li, P.; Pan, X.; Wang, M.; Shao, Y.; Mo, Y.; Mei, Y. Perspective: Reference-Potential Methods for the Study of Thermodynamic Properties in Chemical Processes: Theory, Applications, and Pitfalls. *J. Phys. Chem. Lett.* **2023**, *14*, 4866–4875.