

Automated Discovery and Refinement of Reactive Molecular Dynamics Pathways

Lee-Ping Wang,^{†,‡} Robert T. McGibbon,[†] Vijay S. Pande,[†] and

Todd J. Martinez^{*,†,‡,¶}

[†] *Department of Chemistry, Stanford University, Stanford, CA 94305.*

[‡] *The PULSE Institute, Stanford University, Stanford, CA 94305*

[¶] *SLAC Linear Accelerator Laboratory, Menlo Park, CA 94025*

E-mail: todd.martinez@stanford.edu

November 28, 2015

Abstract

We describe a flexible and broadly applicable energy refinement method, “nebterpolation,” for identifying and characterizing the reaction events in a molecular dynamics (MD) simulation. The new method is applicable to *ab initio* simulations with hundreds of atoms containing complex and multi-molecular reaction events. A key aspect of nebterpolation is smoothing of the reactive MD trajectory in internal coordinates to initiate the search for the reaction path on the potential energy surface. We apply nebterpolation to analyze the reaction events in an *ab initio* nanoreactor simulation that discovers new molecules and mechanisms, including a C-C coupling pathway for glycolaldehyde synthesis. We find that the new method, which incorporates information from the MD trajectory that connects reactants with products, produces a dramatically distinct set of minimum energy paths compared to existing approaches that start from information for the reaction endpoints alone. The energy refinement method

16 described here represents a key component of an emerging simulation paradigm where
17 molecular dynamics simulations are applied to discover the possible reaction mecha-
18 nisms.

19 **1 Introduction**

20 Chemical reactions occur when atoms move through space and rearrange their chemical
21 bonds, and the study of thermodynamics (i.e. reaction energies) and kinetics (reaction
22 rates) of elementary reactions is at the heart of experimental and theoretical chemistry. The
23 essential features of many reactions can be understood by studying the adiabatic potential
24 energy surface (PES) of the electronic ground state, a parametric function of the atomic
25 positions in the high-dimensional configuration space in which the reactant, product, and
26 transition state (TS) structures are stationary points. The transition state is the highest
27 energy point on the minimum energy path (MEP) connecting reactants and products, char-
28 acterized by having exactly one negative eigenvalue in the Hessian matrix. The TS and
29 MEP provide a starting point for understanding the reaction rate of elementary reactions
30 using rate theories such as transition state theory.^{1,2} Numerical methods such as transition
31 path sampling³ and transition-path theory⁴ can incorporate the influence of the thermody-
32 namic ensemble, dynamical effects,⁵ and nuclear quantum effects on the reaction rate. In
33 many mechanistic problems of interest, the reactants and product structures may be known
34 from experiment, but the transient intermediates are relatively difficult to detect and the
35 transition state (TS) structures can only be found by supporting experimental data with the-
36 oretical calculations.^{6,7} This sets the stage for theory and simulation to find and characterize
37 the intermediates and transition states.

38 The search for TS structures is an important challenge in theoretical and computational
39 chemistry that has motivated the development of diverse theoretical methods. The first
40 practical calculations of polyatomic TS structures by Komornicki and coworkers^{8,9} were
41 carried out by first using constrained energy minimizations along manually chosen reaction

42 coordinates to locate an initial guess structure, then minimizing the Euclidean norm of the
43 gradient to locate the TS. However, the gradient is zero at any stationary point on the PES
44 and not just the desired TS, which makes the final answer highly dependent on the initial
45 guess. This difficult problem has motivated the development of methods for generating
46 accurate initial guesses as well as improved optimization algorithms for locating the TS.

47 The procedure for generating an initial guess from known reactant and product struc-
48 tures often begins by constructing a reaction pathway that connects reactant and product
49 (the *endpoints*) using purely geometric methods such as the linear synchronous transit¹⁰ or
50 related interpolation methods.¹¹⁻¹³ Following this, the pathway is iteratively optimized using
51 information from the PES,¹⁴ either by minimizing an energy functional of the pathway as in
52 the nudged elastic band method,^{15,16} or by minimizing the normal component of the force
53 along equally spaced structures as in the string method.^{17,18} The growing string,¹⁸⁻²² freez-
54 ing string,²³ searching string²⁴ and quadratic-PES²⁵ are methods that build and optimize
55 the reaction pathway simultaneously, leading to improved efficiency. We shall refer to these
56 methods collectively as *reaction path-finding methods* that provide a sequence of structures
57 connecting reactant and product along an approximate MEP as well as an initial guess for
58 the TS structure. The TS guess can then be used to initiate a numerical search for the
59 stationary point, using *transition state optimization methods* such as partitioned rational
60 function optimization,²⁶⁻²⁹ the dimer method^{30,31} and other types of minimum mode follow-
61 ing algorithms.³²⁻³⁵ The computational efficiency may be further improved by using internal
62 coordinate systems.³⁶⁻³⁸ Once the TS structure is found, the intrinsic reaction coordinate
63 (IRC) method enables the calculation of the MEP by following the energy downhill (forward
64 and backward) along the mass-weighted steepest descent direction.³⁹⁻⁴³ This is useful to
65 verify that the resultant TS indeed connects the reactant and product structures.

66 Generally speaking, reactivity may be characterized by multiple elementary steps where
67 the intermediates are not known *a priori*, requiring a broader exploration of the configuration
68 space to find new minima on the PES and the pathways connecting them. Global optimiza-

69 tion methods for systematically finding low-energy configurations and the pathways between
70 them include the basin-hopping method by Wales and coworkers,^{44,45} the scaled hypersphere
71 search / global reaction route mapping method by Maeda, Ohno and coworkers,⁴⁶⁻⁴⁹ and
72 the recent minima hopping method.⁵⁰ Reaction discovery may be aided by chemical knowl-
73 edge by applying heuristic rules (e.g. bond breaking) to generate the intermediates^{51,52} or
74 by introducing forces to press reactant and product together.^{53,54} Reaction events can also
75 be observed in *ab initio* molecular dynamics (AIMD) simulations, but a major challenge is
76 the simulation time required to escape from deep free energy minima.⁵⁵ The frequency of
77 reaction events can be greatly accelerated by applying bias potentials that push the sys-
78 tem away from the free energy minima along a collective variable,⁵⁶⁻⁶⁰ which assumes some
79 knowledge of the reaction coordinate or collective variable along which to apply the biasing
80 potential. Reactivity is similarly enhanced when studying extremely high-temperature^{61,62}
81 or high-pressure regimes^{63,64} because that shifts the chemical equilibrium towards products
82 with higher entropy or lower volume, respectively.

83 We recently introduced the *ab initio* nanoreactor, a special type of AIMD simulation
84 focused on exploring and discovering new reaction pathways.⁶⁵ The nanoreactor is a high-
85 temperature AIMD simulation with many reactant molecules in which reactivity is accel-
86 erated by several means: (a) accelerating the electronic structure calculation on graphics
87 processing units,^{66,67} (b) employing an approximate level of electronic structure theory, and
88 (c) a virtual piston which pushes molecules toward the center of the simulation, greatly in-
89 creasing the frequency of collisions and barrier crossings. The nanoreactor MD simulation
90 trajectory contains a great number of *reaction events* that lead from the starting reactants
91 to new and unexpected intermediates and products. This motivates a detailed study of the
92 reaction events to assess their mechanistic relevance for different experimental conditions,
93 including but not limited to high temperature and pressure regimes. Intuitively, a reaction
94 event in the trajectory is a promising initial guess to search for an underlying minimum
95 energy path; however, this calculation is highly challenging because the trajectory usually

96 contains myriad high-frequency and large-amplitude motions that are either orthogonal to
97 the reaction mechanism or involve repeated crossing of the activation barrier. In order to
98 address these challenges, we developed an energy refinement method to calculate minimum
99 energy paths from the nanoreactor MD trajectory. This method, which we call “nebter-
100 polation,” recognizes and extracts reaction events from the simulation trajectory, applies a
101 new internal-coordinate smoothing algorithm to remove high-frequency motions, then ap-
102 plies established reaction path-finding methods and transition state optimizations to afford
103 the transition state and minimum energy path. In this paper, we describe the development
104 of the energy refinement method and its application to several novel examples as well as a
105 large data set of reaction events from a nanoreactor MD simulation.

106 We start by defining and describing the various types of pathways in the stages of energy
107 refinement, starting from the *reaction event* as observed in the nanoreactor, and progress-
108 ing through the *initial pathway* with energy-minimized endpoints, the *smoothed pathway*
109 generated by internal-coordinate smoothing, and the *final pathway* (also the IRC or MEP)
110 from locating the transition state and reconnecting it with the endpoints (Figure 2). Next,
111 we present some interesting representative examples of reactions from the nanoreactor, in-
112 cluding a C-C coupling pathway yielding glycolaldehyde (Figure 3) and some examples of
113 multiple pathways connecting the same reactant and product (Figures 5 and 6). Finally,
114 we investigate the effect of including information from the MD trajectory to find the TS by
115 comparing the outcomes of a *path-based approach* that uses the smoothed pathway to initial-
116 ize the string method,¹⁸ and an *endpoint-only approach* that uses the freezing string method
117 to build the pathway from just the endpoints.²³ We find that both approaches are able to
118 find transition states for a substantial fraction of reactions but with limited overlap, which
119 indicates that using the pathway to supply the initial guess can provide distinct mechanistic
120 information compared to starting from only the endpoints.

2 Methods

2.1 Identifying reaction events

The input to our analysis procedure is an *ab initio* molecular dynamics (MD) trajectory, consisting of a discrete time series of atomic coordinates separated in time by a fixed sampling interval (our coordinates are saved every time step, or 0.5 fs). In principle, any reactive MD method could be used to generate these trajectories, such as reactive force fields⁶⁸ or tight-binding density functional theory⁶⁹ (DFTB). The starting structure in the simulation consists of many different molecules with more than 100 total atoms. The trajectory contains many individual reaction events, but we start with no information describing when the reactions occur or which atoms participate. We shall assume in the following that the MD trajectory contains individual *reaction events* involving a small subset of atoms and taking place in a small time interval; thus, the first task is to identify and extract these reaction events from the simulation to characterize them in greater detail.

We represent molecules using a connectivity graph data structure where nodes represent atoms and edges represent covalent bonds. For each frame \mathbf{Q}_t in the MD simulation (where time is measured in units of the sampling interval), we construct a connectivity graph G_t by assigning bonds to pairs of atoms (i, j) separated by distance $r_{ij; t} < 1.4(R_i + R_j)$, where R_i is the covalent radius⁷⁰ of atom i and the factor of 1.4 helps to compensate for transient bond stretching during the MD simulation. The graphs describing individual molecules are obtained by separating the overall graph into its connected component subgraphs. The nodes and edges also contain attributes of the corresponding atoms and bonds (e.g. node attributes include the chemical element and the global index of the atom in the simulation.)

Since we are interested in analyzing trajectories with a large number of reaction events, it is useful to define the following binary time series for each molecular graph m within the

145 overall simulation:

$$E_t^m = \begin{cases} 1 & \text{if } m \in G_t \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

146 where the $m \in G$ notation signifies that m is isomorphic to one of the connected component
147 subgraphs of G . In the test for isomorphism, the nodes are labeled by the atom indices
148 and edges are not labeled; thus, geometric isomers and stereoisomers are considered to
149 be isomorphic and reaction events that preserve the atomic connectivity (e.g. cis-trans
150 isomerization) are not detected. This is done for convenience in the present work and is not
151 a fundamental limitation - it is also possible to further test graphs in order to distinguish
152 between isomers. The existence of a molecule m in the simulation at time t is thus measured
153 by the binary state of E_t^m .

154 In practice, the MD simulation contains large-amplitude molecular vibrations and col-
155 lisions that briefly perturb the connectivity graphs without any reactivity taking place,
156 manifesting as noise in E_t (we have omitted the superscript m for brevity). We addressed
157 this problem using a two-state hidden Markov model (HMM), in which our observed time
158 series E_t is modeled as a probabilistic function of a two-state Markov chain X . The Markov
159 chain is a stochastic binary sequence of states X_t characterized by a *transition probability*
160 $\mathbf{T} \equiv P(X_t|X_{t-1})$ describing the frequency of transitions in the sequence.⁷¹ Although the
161 states X_t are not directly observable, they generate the observed values of E_t according to
162 the *output probability* $\mathbf{O} \equiv P(E_t|X_t)$ describing how often a given value of X_t produces a
163 particular observation of E_t . The HMM is thus parameterized by the transition probability,
164 the output probability, and the initial probability $P(X_0)$. The end goal is to calculate the
165 most likely sequence of states $\{V_0, V_1, \dots, V_t\}$ through the Markov chain, which contains a
166 reduced amount of noise due to the parameterization of the HMM.

167 We set the initial probabilities to a uniform distribution (i.e. $P(X_0 = 0) = P(X_0 = 1) =$

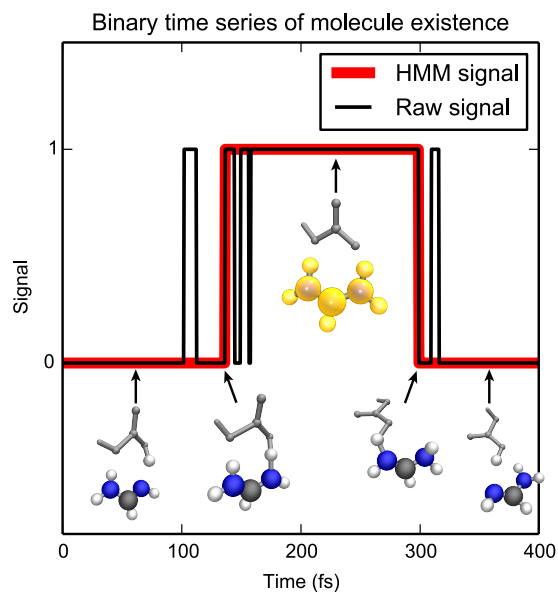


Figure 1: Using a two-state HMM to model the time series of the existence of a molecule. The raw signal (black line) is calculated by building a connectivity graph of the atoms as a function of time and testing for the subgraph corresponding to the molecule highlighted in yellow ($\text{CH}(\text{NH}_2)_2^+$, center). High-amplitude vibrations in the MD trajectory introduce significant noise into the raw signal, and the noise is removed by constructing a two-state HMM and estimating the most likely sequence of hidden states (red line).

168 0.5) and parameterized the transition probability as:

$$\mathbf{T} = \begin{pmatrix} 0.999 & 0.001 \\ 0.001 & 0.999 \end{pmatrix} \quad (2)$$

169 which indicates that X_t has a 0.001 probability of making a $0 \rightarrow 1$ or $1 \rightarrow 0$ transition in
170 any frame. Choosing smaller values in the off-diagonal of \mathbf{T} increases the strength of the
171 noise filter and forces X_t to have fewer transitions. Here, the parameter is chosen so that
172 the frequency of transitions is roughly consistent with the piston interval of 4000 frames in
173 the nanoreactor simulation. We similarly parameterized the output probability as:

$$\mathbf{O} = \begin{pmatrix} 0.6 & 0.4 \\ 0.4 & 0.6 \end{pmatrix} \quad (3)$$

174 which signifies that E_t has a 0.6 probability of being equal to X_t at any given time. Choosing
175 larger values in the off-diagonal of \mathbf{O} allows the Markov process to deviate more often from
176 the observed signal and behave according to its intrinsic transition probability. The HMM
177 provides the joint probability distribution over the observed and hidden variables as:

$$P(X_{0:t}, E_{1:t}) = P(X_0) \prod_{i=1}^t P(X_i | X_{i-1}) P(E_i | X_i). \quad (4)$$

178 where $E_{1:t} \equiv \{E_1, \dots, E_t\}$ represents the inclusive sequence of observed values of E from
179 the initial value to any time t . We apply the Viterbi algorithm to compute the most likely
180 sequence of states over the Markov chain, namely:

$$V_{0:T} = \max_{X_0, X_1, \dots, X_T} P(X_{0:T} | E_{1:T}) \quad (5)$$

181 where T is the length of the whole MD trajectory. As shown in Figure 1, V_t^m accurately
182 models the existence of molecule m at time t and removes the noise from the distance-based

183 measurement. The transitions of V_t^m are effective indicators of reaction events involving the
184 molecule m .

185 A reaction event includes the complete set of all molecules that participate in rearranging
186 their bonds. In order to extract a reaction event involving molecule m , we start at a $1 \rightarrow 0$
187 transition of V^m at t_{event} and trace its atoms a_m forward in time to t_{after} , where they have
188 converted to one or more product molecules $\{n; V_{t_{\text{after}}}^n = 1\}$. If the molecular graphs $\{n\}$
189 contain more atoms than a_m , it means that molecule m does not contain all of the product
190 atoms, and multiple reactant molecules must have been involved. The set of atoms involved
191 in the reaction is then expanded to a_n containing all atoms in $\{n\}$, and the reactant molecules
192 are found by tracing the atoms a_n backward in time to t_{before} , prior to the detected reaction
193 event. This process is iterated back and forth until complete and consistent sets of reactant
194 and product molecules are found.

195 The resulting reaction event \mathbf{q}_t contains the coordinates for the reactant and product
196 atoms over the time interval of the reaction ($t_{\text{before}}, t_{\text{after}}$). \mathbf{q}_t is extracted from \mathbf{Q}_t for the
197 subsequent refinement calculations, assuming that only the electrons of A_m are directly in-
198 volved in the reaction. The net charge and spin polarization are approximated by averaging
199 the Mulliken charge and spin populations over the time interval, summing over the atoms,
200 and rounding to the nearest integer. This is a good approximation for most of the organic
201 reactions considered in this paper, and future work will extend this method to include elec-
202 tron donors and acceptors that participate in the reaction without undergoing changes in
203 their connectivity.

204 Since \mathbf{q}_t contains far fewer atoms and frames than the whole trajectory \mathbf{Q}_t , we may carry
205 out the refinement calculations at a higher level of electronic structure theory than was
206 used to run the MD simulation. In this paper, the MD simulations were carried out using
207 Hartree-Fock (both restricted and unrestricted) and the 3-21G basis set, while the refinement
208 calculations were carried out using unrestricted B3LYP and the 6-31+G(d,p) basis set. Even
209 higher level electronic structure methods could be employed if desired, but this suffices for

210 illustration of the method.

211 2.2 Determining path endpoints

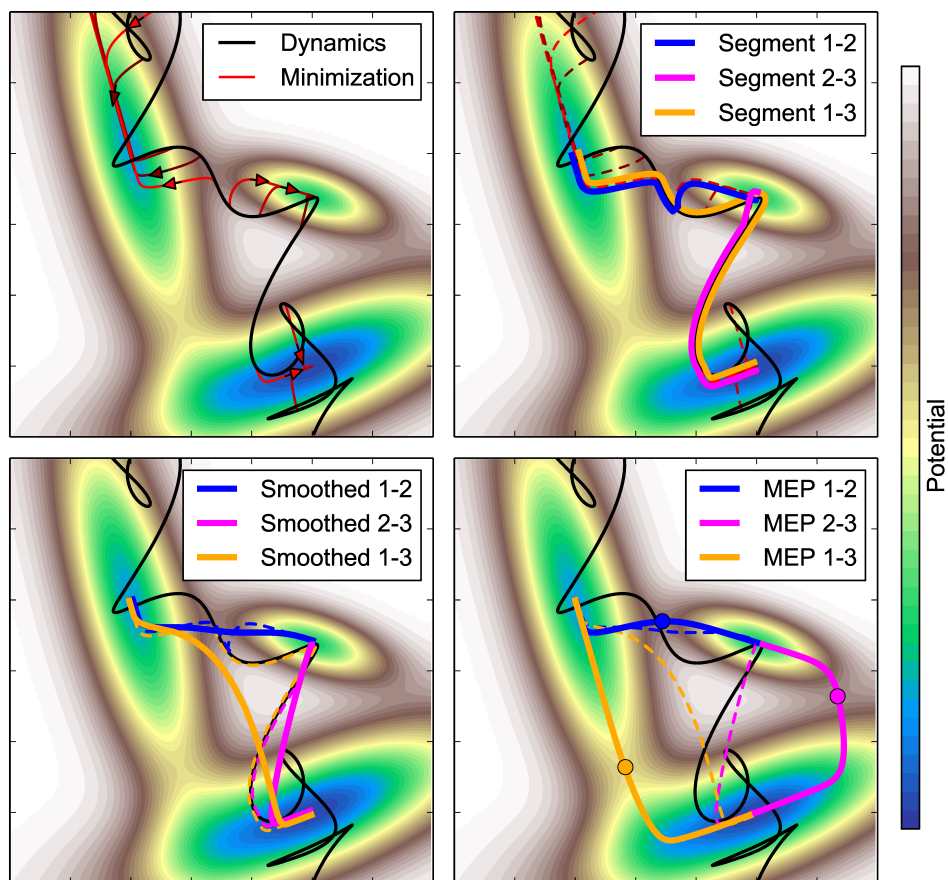


Figure 2: Illustration of the reaction path refinement method. *Top left*: A molecular dynamics trajectory (black curve) traverses the potential energy surface and passes through several energy basins. After identifying this trajectory as a reaction event, evenly sampled points on the trajectory are energy-minimized (red lines with arrows) in order to identify the distinct basins (chemical species). *Top right*: The sequence of structures from the minimization are concatenated with intervals of the molecular dynamics trajectory to form initial pathways connecting the energy basins. *Bottom left*: The connecting segments are smoothed in internal coordinates to decrease their total arc length and curvature. *Bottom right*: Reaction path-finding methods such as nudged elastic band and the string method may be applied to find the transition state, followed by a transition state optimization and intrinsic reaction coordinate calculation to obtain the final MEP.

212 The next step after identifying and extracting the reaction event is to search for energy-
213 minimized structures which will serve as initial guesses to the endpoints of the minimum

214 energy path. This is done by initializing energy minimizations from evenly sampled time
215 points in \mathbf{q}_t ; we used a 10 frame sampling interval (5 fs simulation time), which is rather
216 fine-grained and slightly shorter than an O-H classical vibrational period. Each energy
217 minimization initialized from \mathbf{q}_t provides a sequence of structures that proceeds downhill
218 on the potential energy surface in small steps (on the order of 0.3 Å) and ends at a local
219 minimum. The set of configurations used to initiate the minimizations typically converge
220 to a much smaller set of distinct local minima, allowing us to partition the time series into
221 a small number of intervals - namely, $(t_{\text{begin}}^a, t_{\text{end}}^a)$ for all frames that minimize to chemical
222 species a . An individual reaction event may contain two or more intervals corresponding to
223 reactants, products and intermediates, or only one interval if the minimizations all converge
224 to a single chemical species; the latter may be due to differences in the level of theory between
225 the MD and refinement calculations, the removal of noncovalent interactions when extracting
226 the reactant and product atoms from their environment, or the absence of an energy barrier
227 as is often the case for homolytic bond dissociation to two radicals. In this work, we will
228 focus on reaction events that have at least two energy basins.

229 Because two intervals are likely to be separated by a barrier crossing, we select the
230 sequence of frames in between $(t_{\text{end}}^a, t_{\text{begin}}^b)$ and concatenate it with the two sequences of
231 structures from their associated energy minimizations. This provides a closely spaced se-
232 quence of structures connecting the chemically distinct energy basins (Figure 2, top right),
233 which we call the *initial pathway* \mathbf{x}^{ab} . The initial pathway contains closely spaced frames
234 taken from MD simulations and energy minimizations, which may contain useful information
235 for locating the TS.

236 **2.3 Internal coordinate smoothing**

237 The initial pathway is kinked at the concatenation points and contains high-frequency mo-
238 tions from the MD simulation that render it unsuitable for direct initiation of a minimum
239 energy path search. In this section, we describe a method for smoothing the coordinates

240 in redundant internal coordinates that preserves the energy-minimized endpoints and essen-
241 tial features of the connecting pathway while removing kinks and high-frequency motions.
242 The goal is to transform the initial pathway into a *smoothed pathway* that serves as a good
243 initial guess for reaction path-finding methods. In principle, the Cartesian coordinates can
244 be smoothed directly by simply taking a running average, but since the potential energy
245 surface is highly nonlinear and anharmonic, a linear smoothing could generate unphysical
246 structures with very high energies (e.g. a linear interpolation may result in unphysical struc-
247 tures where atoms pass through each other.) This motivates the use of internal coordinates
248 (e.g. interatomic distances, angles, and dihedrals) for the smoothing procedure.

249 We employed a redundant internal coordinate system in the smoothing procedure, be-
250 cause nonredundant internal coordinate systems are well-known to exhibit singularities in
251 the Jacobian — especially when large changes in the Cartesian coordinates are involved.
252 Our choice of redundant internal coordinates is the union of: (1) all pairwise interatomic
253 distances, and (2) the interatomic distances, angles, and dihedral angles from the union of all
254 bonds that occur in the initial pathway. Since this is an overdetermined coordinate system,
255 the smoothed internal coordinates cannot be uniquely inverted to obtain a set of Cartesian
256 coordinates. We thus define our *inverse Cartesian coordinates* in terms of a least-squares
257 objective function, where the minimum solution is a set of Cartesian coordinates that closely
258 corresponds to the smoothed internal coordinates in a least-squares sense.

259 Our smoothing procedure follows these steps:

260 1. Prior to smoothing, the structures along the initial pathway are linearly spaced along
261 the total arc length measured using the RMS displacement between adjacent frames.

262 The following procedure is applied to each frame in sequence:

263 2. Calculate the redundant internal coordinates $\mathbf{z}_i \equiv IC(\mathbf{x}_i)$ for each structure along the
264 initial pathway.

265 3. Calculate smoothed redundant internal coordinates for each structure $\tilde{\mathbf{z}}_i$ by convolution

266 with windowing function: $\tilde{\mathbf{z}}_i = \sum_{j=-\Delta}^{+\Delta} \mathbf{z}_{i+j} w(j)$, where the windowing function $W(j)$ is
 267 defined on the interval $(-\Delta, +\Delta)$ and normalized to one.

268 4. Find the inverse Cartesian coordinates that minimize the least-squares error for each
 269 set of smoothed redundant internal coordinates: $\tilde{\mathbf{x}}_i = \min_{\mathbf{x}'} (IC(\mathbf{x}') - \tilde{\mathbf{z}}_i)^2$, using \mathbf{x}_i as
 270 the initial guess.

271 5. The sequence of inverse Cartesian coordinates is the smoothed pathway.

272 We found that introducing a repulsive pseudo-energy term in the objective function was
 273 helpful for preventing close contacts between pairs of atoms during the minimization. With
 274 the repulsive term, the total function to be minimized takes the following form:

$$\begin{aligned} \chi^2(\mathbf{x}) = & \sum_{i,j \in N} |w_r(r_{ij}(\mathbf{x}) - \tilde{r}_{ij})|^2 + w_V V_{ij}(r_{ij}(\mathbf{x})) \\ & + \sum_{\{i,j,k\} \in \text{angles}} \left| w_\theta(\theta_{ijk}(\mathbf{x}) - \tilde{\theta}_{ijk}) \right|^2 + \sum_{\{i,j,k,l\} \in \text{torsions}} \left| w_\phi(\phi_{ijkl}(\mathbf{x}) - \tilde{\phi}_{ijkl}) \right|^2, \end{aligned} \quad (6)$$

275 where

$$V_{ij}(\mathbf{x}) = \begin{cases} D_{ij}(1 - e^{-a_{ij}(r_{ij} - r_{ij}^0)})^2 & : x < r_{ij}^0 \\ 0 & : x \geq r_{ij}^0. \end{cases} \quad (7)$$

276 Here, $r_{ij}(\mathbf{x})/\theta_{ijk}(\mathbf{x})/\phi_{ijkl}(\mathbf{x})$ are interatomic distances/angles/dihedrals (respectively) cal-
 277 culated from the trial coordinates \mathbf{x} , and $\tilde{r}_{ij}/\tilde{\theta}_{ijk}/\tilde{\phi}_{ijkl}$ are the smoothed internal coordinates
 278 which are the target values in the minimization. V_{ij} is the repulsive part of the Morse poten-
 279 tial which goes smoothly to zero at r_{ij}^0 , and the parameters D_{ij} , a_{ij} and r_{ij}^0 are taken from
 280 standard tables of bond dissociation energies, bond lengths and vibrational force constants.
 281 The weights $w_r = 1.0 \text{ \AA}^{-1}$, $w_\theta = (\frac{\pi}{6})^{-1}$, $w_\phi = \pi^{-1}$ and $w_V = (0.01 \text{ kJ/mol})^{-1}$ are chosen such
 282 that all contributions to the objective function are the same order of magnitude in numerical
 283 tests.

284 The independent minimizations for each $\tilde{\mathbf{x}}_i$ are nonlinear and the objective function
 285 depends on the smoothed internal coordinates $\tilde{\mathbf{z}}_i$; thus, it is possible to find adjacent pairs
 286 of inverse Cartesian coordinates $(\tilde{\mathbf{x}}_{i-1}, \tilde{\mathbf{x}}_i)$ with large deviations, resulting in a undesirable
 287 discontinuity in the smoothed pathway. When this happens, we avoid discontinuities by
 288 introducing a restraint term into the minimization:

$$\bar{\chi}_i^2(\mathbf{x}) = \chi_i^2(\mathbf{x}) + \bar{w}|\mathbf{x} - \tilde{\mathbf{x}}_{i-1}|^2, \quad (8)$$

where $\bar{\chi}^2$ is the objective function with the restraint, and \bar{w} is the weight associated with
 the restraint that biases the solution towards the previous frame. The minimizations are
 carried out sequentially so that $\tilde{\mathbf{x}}_{i-1}$ is known when starting the search for $\tilde{\mathbf{x}}_i$. Starting with
 a restraint term of zero, a discontinuity is detected using the following condition:

$$\frac{\max \text{abs}(\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_{i-1})}{\max \text{abs}(\mathbf{x}_{i+1} - \mathbf{x}_i)} > 2$$

289 , and the search for $\tilde{\mathbf{x}}_i$ is restarted with a small weight factor $\bar{w} = 0.02$. The weight \bar{w} is
 290 increased by successive factors of 1.5 until the discontinuity is suppressed.

291 Introducing the restraint term can lead to problematic behavior in a very small number of
 292 cases where the ending structure on the smoothed pathway is different from that of the initial
 293 pathway. When this happens, we recalculate a sequence of inverse Cartesian coordinates
 294 starting from the ending structure and going in reverse. The search is finished when a
 295 structure from the reversed pathway coincides with a structure on the forward pathway. We
 296 found that this method was robust in smoothing a diverse set of initial pathways where linear
 297 interpolation gives unphysical interatomic distances of less than 0.5 Å.

298 The *smoothed pathway* consists of a set of Cartesian coordinates $\{\tilde{x}_1, \tilde{x}_i, \dots, \tilde{x}_n\}$ which we
 299 use to initialize the minimum energy path search described in the next section. The code
 300 for carrying out the internal coordinate smoothing operation, *nebterpolator*, is open-source
 301 and freely available on the Web.

302 2.4 Searching for the minimum energy path

303 The fourth and final step is to locate the minimum energy path (MEP). We employ standard
304 methods, such as nudged elastic band (NEB) or the string method, using the smoothed
305 pathway as an initial guess. These methods iteratively search for a discretized pathway that
306 minimizes an energy functional of the path (in the case of NEB) or the component of the
307 gradient perpendicular to the path (in the case of the string method). The output of path
308 optimization is a discretized pathway that approximates the MEP and a transition state (TS)
309 estimate corresponding to the highest energy structure along the pathway. The TS estimate
310 is then used to initialize a TS geometry optimization, which searches for a critical point
311 with exactly one imaginary mode, and the character of the TS is verified using a frequency
312 calculation. Finally, the optimized TS geometry connects back to the reactant and product
313 structures via the imaginary mode that leads downhill in both directions; a pair of intrinsic
314 reaction coordinate (IRC) calculations provide the final pathway connecting the TS to the
315 reactant and product energy basins.

316 Our proposed approach uses the smoothed pathway to initiate the search for a minimum
317 energy path; this contrasts with methods such as growing string and freezing string that
318 proceed from knowledge of only the endpoints of the initial pathway. Since our approach
319 provides the endpoints as well as the connecting pathway, we compared our *pathway-based*
320 *approach* (using the smoothed pathway to initiate a string method search for the TS) with
321 an *endpoint-only approach* (using the endpoints of the initial pathway to initiate a freezing
322 string search for the TS). Our results compare the pathway-based approach and the endpoint-
323 only approach to assess whether the connecting pathway provides additional information for
324 locating the TS.

3 Computational Details

The calculations in this paper fall into two categories. In the first category, the nanoreactor AIMD simulations generate trajectories with hundreds of atoms and a large number of reaction events; these simulations were carried out using the TeraChem quantum chemistry software package.⁶⁷ In the second category, the refinement calculations include smaller-scale quantum chemistry calculations such as geometry optimizations, reaction path-finding and transition state searches carried out using the Q-Chem quantum chemistry software package.⁷² In particular, the freezing string calculations in the endpoint-only approach were performed in a delocalized internal system with 21 nodes, 3 gradients per step, LST interpolation and a quasi-Newton optimizer, following the suggested settings from the Q-Chem documentation. The refinement calculations also involve a large amount of geometric analysis and manipulations, as well as automation for carrying out the calculations on more than a thousand reaction events. The automation was implemented in a Python module that communicates with Q-Chem and also includes *nebterpolator*. The parallel refinement calculations used the Work Queue distributed computing framework (Figure S2).

Table 1: Summary information of nanoreactor simulations analyzed in this work. The parameters of the boundary potential are described in the main text. Also listed are the number of reaction events found in each simulation and the corresponding number of final pathways; each reaction event can lead to any number of final pathways.

Label	Length (ps)	Method	Basis	Charge	t_1 / t_2 (ps)	r_1 / r_2 (Å)	k_1 / k_2 (kcal mol ⁻¹ Å ⁻² amu ⁻¹)	Reaction Events	Final Paths
1	59.69	RHF	STO-3G	0	1.0 / 1.0	14.4 / 7.2	1.0 / 1.0	62	90
2	17.94	RHF	STO-3G	0	2.0 / 0.5	14.4 / 7.2	1.0 / 1.0	29	19
3	15.24	RHF	3-21G	0	2.0 / 0.5	14.4 / 7.2	1.0 / 1.0	19	3
4	253.40	RHF	STO-3G	0	1.5 / 0.5	14.0 / 8.0	1.0 / 0.5	47	22
5	45.94	UHF	STO-3G	-1	1.5 / 0.5	14.0 / 9.0	1.0 / 0.5	31	64
6	78.17	UHF	STO-3G	-1	3.5 / 0.5	14.0 / 9.0	1.0 / 0.5	35	7
7	436.16	RHF	3-21G	0	1.5 / 0.5	14.0 / 8.0	1.0 / 0.5	317	713
8	277.00	UHF	3-21G	-1	1.5 / 0.5	14.0 / 9.0	1.0 / 0.5	127	47
9	37.11	UHF	3-21G	0	1.5 / 0.5	14.0 / 8.0	1.0 / 0.5	53	62
10	17.94	UHF	STO-3G	0	1.5 / 0.5	14.0 / 8.0	1.0 / 0.5	8	14
11	155.29	RHF	3-21G	0	1.5 / 0.5	14.0 / 8.0	1.0 / 0.5	249	59
12	65.80	RHF	3-21G	0	1.5 / 0.5	14.0 / 9.0	1.0 / 0.5	43	17

340 Table 1 provides the details of the nanoreactor simulations that are analyzed in this
341 paper. The simulations used either the restricted or unrestricted Hartree-Fock (RHF/UHF)
342 electronic wavefunction and small Gaussian basis sets (STO-3G or 3-21G) to calculate the
343 Born-Oppenheimer potential energy surface. Four different initial configurations were used
344 containing the same molecules: 14 H₂O, 14 CH₄, 14 NH₃, 14 CO, 16 H₂. The equations of
345 motion were numerically integrated using Langevin dynamics with a time step of 0.5 fs, an
346 equilibrium temperature of 2000 K (also the starting temperature), and a friction coefficient
347 of 7 ps⁻¹. A flat-bottom spherical potential was used to contain the molecules within a finite
348 volume as $U(r) = mk(r - r_0)^2\theta(r - r_0)$, where r_0 is the sphere radius, k a force constant, m
349 the atomic mass, and $\theta(r - r_0)$ is the Heaviside step function. This potential is zero out to r_0
350 and quadratic for larger values of r . The parameters of $U(r)$ are modulated as a rectangular
351 pulse waveform; the parameters are held at r_1, k_1 for time t_1 , and then "pulsed" to r_2, k_2 for
352 time t_2 ; importantly, the smaller value of r_2 forces molecules with radial distance $r > r_2$ to
353 accelerate towards the center and collide at high velocities, which drives the reactivity. The
354 spherical potential is proportional to the atomic mass to ensure equal inward accelerations
355 of molecules that are outside of the sphere during these compression phases. A total of
356 1020 reaction events were found in 1460 ps of aggregate simulation time; the single longest
357 simulation ran for 436 ps and produced 317 reaction events.

358 The energy refinement calculations used the B3LYP hybrid density functional and two
359 Gaussian basis sets; the 6-31G(d) basis (here referred to as the small basis) was used in all
360 calculations leading up to and including the transition state search. The 6-31+G(d,p) basis
361 (here referred to as the large basis) was used to reoptimize the transition state from the
362 small basis calculation and calculate the IRC leading from the transition state to the nearest
363 energy minimum.

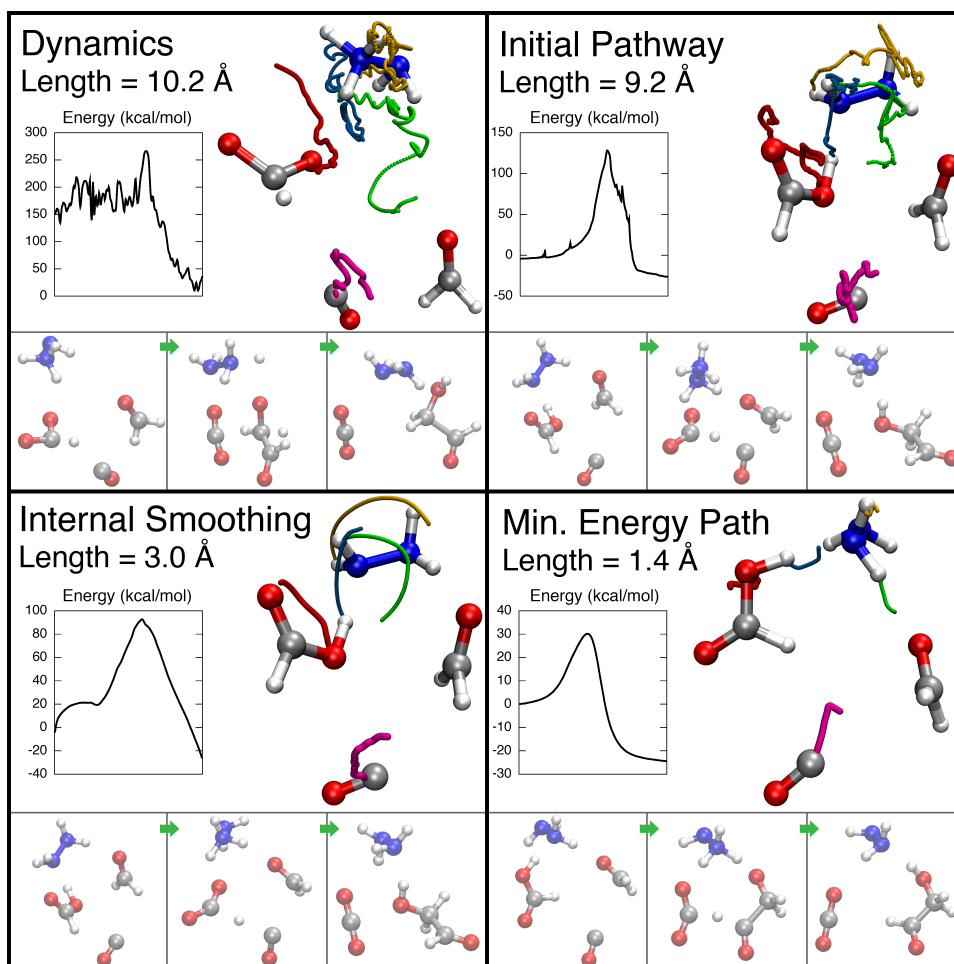
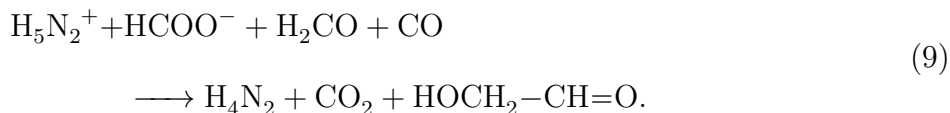


Figure 3: Demonstration of reaction pathway refinement applied to the glycolaldehyde example reaction. Each panel indicates one stage of the refinement. The initial structure of the pathway is shown in ball-and-stick representation (C, gray; O, red; H, white), and colored tracks indicate the positions of selected atoms along the pathway (magenta, C; red, O; gold/green/blue, H). Three frames taken from the start, midpoint, and end of the pathway are rendered at the bottom of each panel. The path length is given by the cumulative sum of the RMS displacement between frames. The inset shows the potential energy calculated at the B3LYP/6-31+G** level along the pathway. *Top left*: Frames containing reaction event selected from the nanoreactor MD simulation. *Top right*: Initial pathway constructed by connecting structures from energy minimization and molecular dynamics. *Bottom left*: Internal coordinate smoothing of the initial pathway. *Bottom right*: Final minimum energy path, optimized using the smoothed pathway as an initial guess. The reaction energy and activation energy are -24.4 kcal/mol and 30.1 kcal/mol.

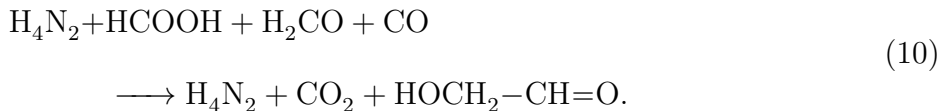
364 4 Results and Discussion

365 4.1 Example reaction: Glycolaldehyde retrosynthesis

366 In this section, we describe an example of the energy refinement procedure applied to a reac-
367 tion event observed in the nanoreactor, where glycolaldehyde ($\text{HOCH}_2\text{-CH=O}$) and carbon
368 dioxide (CO_2) molecules collide to form carbon monoxide (CO), formaldehyde ($\text{H}_2\text{C=O}$) and
369 formate ion (HCOO^-). Hydrazine (H_4N_2) participates in the collision and accepts a proton
370 to form hydrazinium ion (H_5N_2^+). Since the reaction pathway is equally valid in either di-
371 rection, we will treat glycolaldehyde as the product in the subsequent analysis because the
372 discussion is more chemically intuitive. Thus, the (reversed) reaction event is given by:



373 The reaction event is 340 frames (170 fs) in length, with 17 atoms and an overall formula
374 of $\text{C}_3\text{H}_8\text{N}_2\text{O}_4$. The sum of the Mulliken charge populations over the atoms has a mean of
375 0.117 over the frames and a standard deviation of 0.140, and the Mulliken spin populations
376 were 0.0 throughout (the MD simulation used restricted HF). Based on this, we assigned
377 neutral charge and singlet multiplicity for subsequent calculations on this pathway. Energy
378 minimizations of the endpoints cause hydrazinium to transfer a proton to formate, so the
379 reaction event after minimization is:



380 The endpoints of the initial pathway are sufficient to search for the transition state via
381 the endpoint-only approach. The endpoint-only approach did not produce a viable reaction
382 pathway or transition state estimate, as the highest energy along the path was > 10 a.u.
383 higher than the energy of the reactants, and much higher than any frame along the dynamics

384 pathway (< 0.5 a.u. higher than the reactants). The high energies were due to close contacts
385 along the freezing string pathway, as the closest approach between any pair of atoms was
386 0.23 \AA in the calculation output. It should be stressed that the freezing string calculations
387 were carried out on thousands of reactions with the same settings as described in Section 3,
388 indicating the success of the method for many other cases. Similarly, a linear interpolation
389 of the Cartesian coordinates from the initial to the final structure resulted in a close contact
390 of 0.17 \AA (Figure S1). The complex motions of atoms in the reaction event indicate a high
391 degree of difficulty in finding the reaction pathway using the endpoints alone.

392 We next applied the path-based approach; internal-coordinate smoothing led to the path-
393 way in Figure 3, bottom left. The smoothed pathway connects the reactant and product
394 through a curved and much shorter pathway, with a total arc length of 3.0 \AA (vs. 9.2 \AA
395 for the initial pathway.) The energies along the pathway are also lower, with a maximum
396 energy of 92 kcal/mol (vs. 128 kcal/mol for the initial pathway.) The reduction in maximum
397 energy is an encouraging result and shows that internal coordinate smoothing can generate a
398 chemically reasonable connecting pathway without the oscillations and kinks from the initial
399 pathway.

400 The smoothed pathway was used to initiate a string method calculation, which provides
401 a starting structure to the TS optimization. The IRC calculation starting from the TS is
402 shown in Figure 3, bottom right, and provides the final pathway connecting the TS with
403 its reactant and product basins. Examination of the final pathway reveals some interesting
404 comparisons with the earlier stages of refinement. The endpoints of the final pathway are
405 chemically identical to those of the initial / smoothed pathway, and the atomic indices of
406 the transferred hydrogen atoms are the same, yet we observed two major differences:

407 First, the smoothed pathway has a greater arc length than the MEP as it contains a
408 tumbling motion of hydrazine and formic acid. This is likely due to the potential energy sur-
409 faces of the reactant and product basins containing numerous local minima from nonbonding
410 interactions that are well-separated in space but close in energy. In other words, the MEP

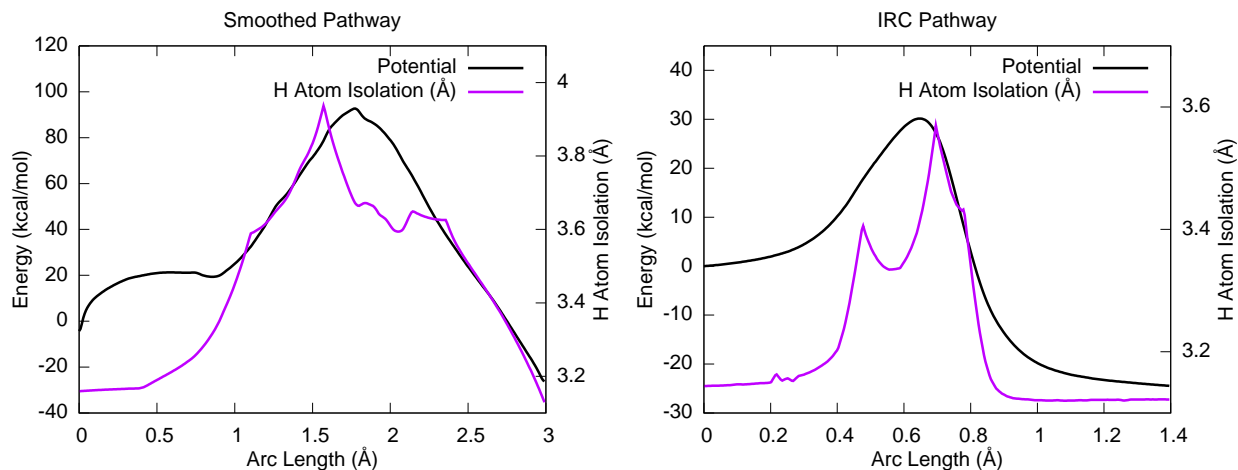


Figure 4: Plot of the potential energy along the smoothed pathway and final pathway from Figure 3, and a simple metric describing the degree of isolation of transferred hydrogen atoms, defined in the main text. The coincidence of the curves in the high-energy region indicates that the highest energies along the smoothed pathway may be due to the occurrence of isolated hydrogen atoms, a kind of chemically unfavorable structure resulting from smoothing.

411 that connects the endpoints of the initial pathway actually crosses several small barriers
 412 corresponding to covalent and noncovalent rearrangements, and the IRC calculation only
 413 finds the two basins nearest the TS. In this example we successfully found the chemically
 414 relevant transition state; we address the frequency of success in Section 4.3.

415 Second, the transition state energy of the final pathway is much lower at just 30 kcal/mol,
 416 compared to 92 kcal/mol for the smoothed pathway. We observed that the structures in the
 417 smoothed pathway contained H atoms that were distant from both the donor and acceptor
 418 site (i.e. r_{DH} and r_{HA} were both large), whereas in the final pathway the donor and acceptor
 419 approached each other prior to H-transfer such that r_{DH} and r_{HA} are never both large. Figure
 420 4 plots the energy on the smoothed and IRC pathways, along with a simple measure of the
 421 hydrogen atom isolation, given by the sum of $\min(r_{\text{DH}}, r_{\text{HA}})$ for each of the three H atoms
 422 being transferred. The qualitative similarity between the hydrogen atom isolation and the
 423 potential in the high-energy regions indicates that the smoothing procedure may benefit
 424 from having even more chemical information, for example from a lower level of QM theory
 425 or a reactive force field. We leave this possibility open for future research.

426 **4.2 Single reaction, multiple pathways**

427 The final pathway found by the refinement procedure is a local MEP in the space of all possi-
 428 ble pathways that connect two minima on the potential energy surface. Since the nanoreactor
 429 simulations often finds several reaction events with the same reactant and product molecules,
 430 it could discover several distinct pathways for the same reaction. In what follows, we discuss
 431 some reaction pathways that connect two minima separated by a single barrier where the
 432 pathways are qualitatively different.

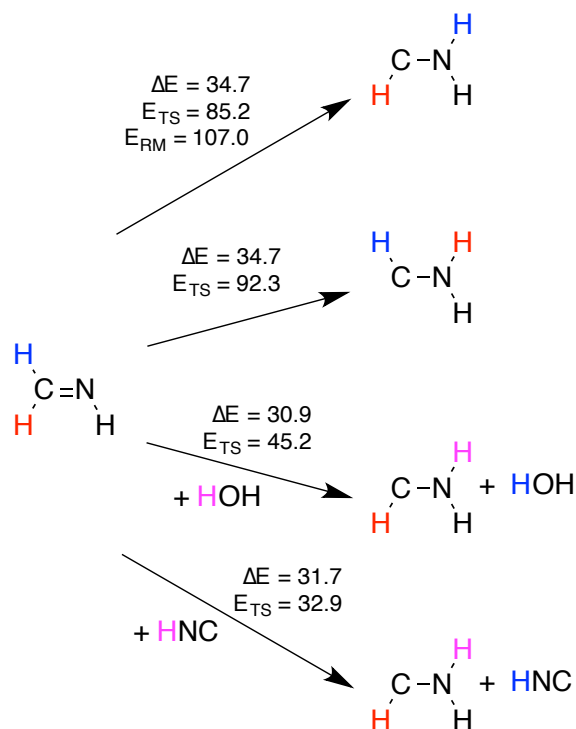


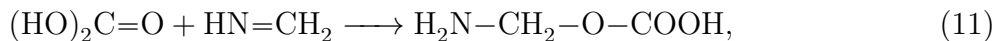
Figure 5: Several isomerizations of formaldimine to aminomethylene. Hydrogen atoms are colored to indicate the pathway that is followed. E_{RM} denotes the plateau energy of the roaming pathway in which H dissociates before attaching to N. The first and second pathway differ by which hydrogen migrates; in the top pathway the initial and final H positions are in a *cis* configuration, and has a lower barrier than the *trans* configuration. In the bottom two pathways, a proton is exchanged with H_2O and HNC.

433 An interesting example is the isomerization of singlet formaldimine to aminomethylene:
 434 $CH_2=NH \longrightarrow CH-NH_2$, involving the migration of a single H atom with a calculated

435 reaction energy of 34.7 kcal/mol. This reaction has several distinct pathways (Figure 5), two
436 of which involve migration of the H in the *cis* or the *trans* position with respect to its final
437 position. The *cis* pathway has a barrier of 85.2 kcal/mol and H moves across in a linear
438 path, whereas the *trans* pathway has a higher barrier of 92.3 kcal/mol and H crosses over
439 the C=N bond in a highly curved path. A “roaming” pathway^{73,74} was also found where
440 the H atom completely dissociates before reattaching; the highest energy is 107.0 kcal/mol
441 above the reactants, and the energy along the pathway resembles a broad plateau rather
442 than a well-defined transition state. The roaming pathway is perhaps better treated as two
443 separate barrier-less reactions involving dissociation and association of H, and we relegate
444 further discussions of such reactions to a future paper.

445 Many reaction pathways are affected by the presence of water, ammonia or other proton-
446 labile molecules that may participate as a bridge in proton / hydrogen atom transfer. In these
447 cases, the pathway involving the participating proton bridge always has a lower barrier in our
448 Urey-Miller nanoreactor data set. When water participates in formalimine isomerization,
449 the barrier is dramatically lower at 45.2 kcal/mol, and the reaction energy changes to 30.9
450 kcal/mol due to noncovalent interactions. When hydrogen isocyanide (HNC) participates,
451 the barrier decreases even further (32.9 kcal/mol) and is almost completely diminished as
452 the reaction energy is 31.7 kcal/mol. By contrast, hydrogen cyanide (HCN) reacts in a very
453 different way by adding to formalimine to form $\text{H}_2\text{N}-\text{CH}_2-\text{N}\equiv\text{C}$, which can then isomerize
454 to form 2-aminoacetonitrile ($\text{H}_2\text{N}-\text{CH}_2-\text{C}\equiv\text{N}$).

455 In most of the cases we examined, multiple pathways connecting the same two endpoints
456 differ by a proton or H-atom bridge. Examples involving differences in heavy atom behavior
457 are less common; here we provide such an example of how heavy atoms may behave differently
458 along a reaction pathway. The nucleophilic addition of carbonic acid to formalimine leads
459 to aminomethyl hydrogen carbonate:



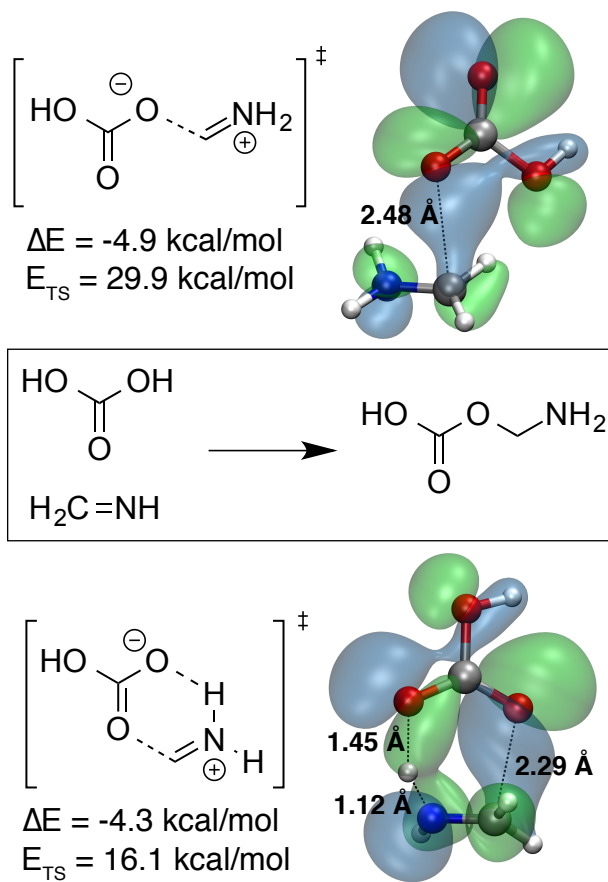


Figure 6: Transition states for two different reaction pathways for the nucleophilic addition of carbonic acid to formaldimine. The upper transition state has an energy that is more than 10 kcal/mol higher than the lower transition state. Differences in reaction energies are due to small conformational differences in the reactant and product molecules. The highest occupied Kohn-Sham orbital is rendered at an isosurface value of ± 0.02 . Key distance measurements are given in Angstrom. Atoms are colored as C, gray; O, red; N, blue; H, white.

460 in which a C–O bond is formed and one hydroxyl proton is transferred to the imine. Two
461 qualitatively different pathways were found with chemically equivalent endpoints (Figure 6):
462 a “lower” pathway with $\Delta E = -4.3$ kcal/mol, $E_a = 16.1$ kcal/mol and an “upper” pathway
463 with $\Delta E = -4.9$ kcal/mol, $E_a = 29.9$ kcal/mol. In the lower pathway, the proton donor and
464 the nucleophile are two different O atoms on carbonic acid, whereas in the upper pathway
465 the same O-atom first transfers the proton then migrates over to make the C–O bond. The
466 difference in barrier heights may be rationalized by looking at the highest occupied Kohn-
467 Sham orbital of the transition states; the HOMO in the upper pathway has C–O bonding
468 character, whereas the HOMO in the lower pathway also has O–H and N–H bonding char-
469 acter. Thus, the transition state in the lower pathway avoids complete dissociation of the
470 O–H bond until the C–O bond is formed.

471 In the preceding reactions with multiple distinct pathways, the products at the end of the
472 pathways are chemically identical but different in terms of the atomic indices; an example is
473 the two isomerization pathways of formaldimine on the top of Figure 5, where two H atoms
474 are swapped. It may be difficult to experimentally measure the relative participation of the
475 pathways, but isotope labeling studies could offer a possible route towards distinguishing
476 them. In any case, it is important to consider all of the pathways when deriving a rate
477 expression for any elementary reaction.

478 4.3 Statistical trends in method behavior

479 We compared the performance and overall behavior of the pathway-based and endpoint-only
480 approaches by performing calculations on a diverse set of systems. We investigated a set
481 of 2652 initial pathways; the results are summarized in Figure 7, where the colored squares
482 indicate the frequencies of the different possible outcomes in each approach. The blue square
483 indicates the number of systems where both approaches found a final pathway containing a
484 reaction, and the orange (resp. green) square counts the systems where only the path-based
485 (resp. endpoint-based) approach produced a positive result. The bottom right quadrant

		Path-based method results:				
		Positive		Negative		
Endpoint-based method results:	Positive	489	77	141		Match
	Negative	26	238	152		Different
	Negative	77		404		No TS
				256	204	
		202	386			
		Match	Different	No TS	No Reaction	

Figure 7: Overall outcome of energy refinement carried out on a data set of 3282 initial pathways from the Urey-Miller nanoreactor simulation, using the path-based approach and the endpoint-only approach. The blue, orange and green boxes respectively count the systems where both approaches, or only one (path-based or endpoint-based) produced a final pathway containing a reaction. The final pathways are further divided into those that recover the same endpoints of the initial pathway, and those that lead to different endpoints. The gray and white squares count negative results where neither approach could find a final pathway containing a reaction, either because a transition state could not be found or because the final pathway contained no reactions.

486 counts the systems where both approaches produced negative results, either due to failing
487 to find a transition state (gray square) or finding a final pathway with no reactivity (white
488 square).

489 Overall, 1604 calculations out of the total 2652 resulted in reactive final pathways, 810
490 of which correctly led back to the reactant and product species of the initial pathway. 1048
491 calculations produced negative results where no reactive final pathway was found. The
492 number of transition states found using the automated procedures is encouraging given the
493 high difficulty in locating transition states. We observed some interesting comparisons; the
494 pathway-based approach was more often able to find a final pathway containing a reaction
495 (1311 vs. 1123 for the endpoint-based method), but conversely the endpoint-based method
496 was more likely to find a final pathway that matches the endpoints of the initial pathway
497 (707 vs. 592 for the path-based method). Our hypothesis is that the path-based method
498 sometimes provides an initial guess that crosses more than one barrier, indicating that more
499 fine-grained division of the pathways may be possible. On the other hand, the endpoint-based
500 method cuts the corners of the MD pathway by connecting the endpoints more directly, in-
501 creasing the likelihood of finding a single barrier but also may fail when the energy landscape
502 is more complex.

503 The collection of final pathways contains some redundancies in the chemical space because
504 many reactions are observed to happen more than once. Figure 8 is a histogram of distinct
505 final pathways, binned by the mean arc length and number of atoms. There are 1157 distinct
506 pathways in all, and the overlap of the two methods is quite small; only 334 distinct pathways
507 are found by both methods, whereas the path-based (resp. endpoint-based) method alone
508 could find 517 (resp. 306) of the pathways. Furthermore, the path-based method produces
509 a greater number of pathways with long arc lengths ($> 6 \text{ \AA}$), and the relative success of
510 endpoint-based methods increases with the number of atoms. The different distributions of
511 pathways found by either method indicates that both are helpful when the goal is to broadly
512 explore the chemical space.

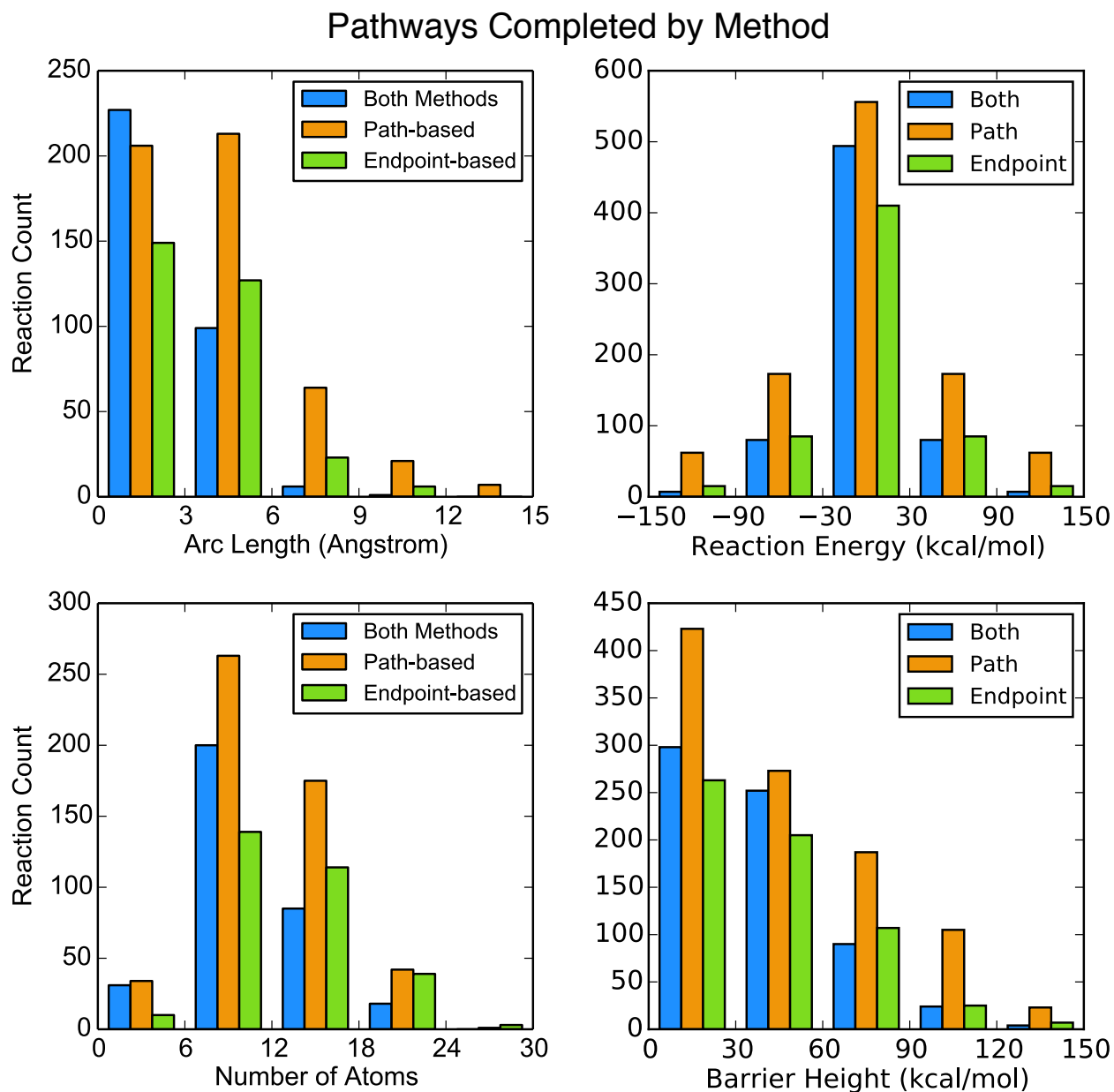


Figure 8: Statistics on distinct final pathways resulting from reaction events in the Urey-Miller nanoreactor simulation. There are a total of 1157 final pathways in the data set with different pairs of endpoints. Green (resp. orange) bars denote IRC pathways that were only found by the endpoint-based method and not the path-based method (resp. vice versa), and blue bars denote IRC pathways that were found by both methods.

5 Conclusion

In the past ten years, *ab initio* molecular dynamics and other kinds of reactive MD simulations have emerged as a promising approach for exploring chemical reactivity in complex systems with many reaction intermediates and elusive transition states. This work describes a systematic approach for connecting the reactive molecular dynamics simulation with automated approaches for mapping individual reaction events in the MD simulation to zero-temperature minimum energy paths. The key components of the method are to recognize the reaction event, locate the energy minima corresponding to the reactant and product, and leverage the MD trajectory to construct a smoothed pathway used to initialize a reaction path-finding calculation. Our calculations indicate that including the trajectory information in a pathway-based approach leads to a greater number of distinct reaction pathways compared to using an endpoint-only approach using only reactant and product structures, indicating that the MD trajectory provides valuable information for finding transition states and reaction pathways in the chemical space.

The method and calculations presented in this paper highlight many promising directions for future work. The application of more detailed electronic structure methods is certainly needed to verify and/or improve the quantitative accuracy of the results. Complete active space approaches^{75–78} are especially relevant when static correlation plays an important role in the reactants or the transition state, but the systematic selection of an active space for such a large number of reactions is a major challenge. Including solvent effects could qualitatively alter the PES and lead to new intermediates and pathways, but the heavy reliance on geometry optimization requires solvent models with a well-behaved potential energy surface; switching-Gaussian implicit solvent models or cluster-continuum models may provide a good first-order description.^{79–81} Hundreds of the initial pathways did not lead to a transition state or MEP containing a reaction; these could serve as a “challenge” data set for further improving reaction path-finding or TS optimization methods. More rigorous techniques for extracting the reaction events from the MD simulation could be helpful for

540 identifying nearby molecules in the simulation that have nontrivial effects on the pathway.

541 Notably, the majority of trajectories used in this paper used restricted HF because we
542 could attain much longer simulation lengths compared to unrestricted HF (UHF). This is
543 mostly due to the relative ease of converging the self-consistent field in RHF, which decreases
544 the computational cost of each time step. The choice of wavefunction affects the reactions
545 discovered; the UHF simulations contained more reaction events involving homolytic bond
546 cleavage and radical species, whereas the RHF simulations were more heterolytic in nature.
547 This is presumably due to the inability of RHF to describe the open-shell electronic state
548 for homolytic bond breaking reactions and overestimating the dissociation energy. Although
549 RHF prevents some reactions from happening, we did not observe a significant decrease in
550 overall reactivity; rather, this seemed to bias the simulation toward finding reactions that
551 did not involve radical or open-shell species, such as nucleophilic attack and proton transfer.
552 These results motivate future studies that investigate how the choices of electronic wavefunc-
553 tion, basis set, and/or DFT approximation could be used to shift the distribution of reaction
554 events for refinement, toward achieving our objective of discovering and characterizing the
555 broad landscape of reactions in complex experimental conditions.

556 **Supporting Information Available**

557 Expanded version of Figure 3, comparing results from different minimum energy path opti-
558 mization methods and flow chart describing distributed algorithm for path refinement. This
559 information is available free of charge via the Internet at <http://pubs.acs.org> This material
560 is available free of charge via the Internet at <http://pubs.acs.org/>.

561 **Acknowledgement**

562 This work was supported by the Office of Naval Research (N00014-14-1-0590) and the DOE
563 Office of Basic Energy Science through the Predictive Theory of Transition Metal Oxide

564 Catalysis Grant. TJM is grateful to the Department of Defense (Office of the Assistant
565 Secretary of Defense for Research and Engineering) for a National Security Science and
566 Engineering Faculty Fellowship (NSSEFF). We are grateful to Matt Harrigan for providing
567 software that assisted in the drawing of Figure 2.

568 **References**

- 569 (1) Berne, B. J.; Borkovec, M.; Straub, J. E. *J. Phys. Chem.* **1988**, *92*, 3711.
- 570 (2) Cukier, R. *J. Phys. Chem. A* **1999**, *103*, 5989.
- 571 (3) Bolhuis, P. G.; Chandler, D.; Dellago, C.; Geissler, P. L. *Ann. Rev. Phys. Chem.* **2002**,
572 *53*, 291.
- 573 (4) E, W.; Vanden-Eijnden, E. *Ann. Rev. Phys. Chem.* **2010**, *61*, 391.
- 574 (5) Ammal, S.; Yamataka, H.; Aida, M.; Dupuis, M. *Science* **2003**, *299*, 1555.
- 575 (6) Williams, I. H. *Chem. Soc. Rev.* **1993**, *22*, 277.
- 576 (7) Ziegler, T.; Autschbach, J. *Chem. Rev.* **2005**, *105*, 2695.
- 577 (8) McIver, J. W.; Komornicki, A. *J. Am. Chem. Soc.* **1972**, *94*, 2625.
- 578 (9) Komornicki, A.; Ishida, K.; Morokuma, K.; Ditchfield, R.; Conrad, M. *Chem. Phys.*
579 *Lett.* **1977**, *45*, 595.
- 580 (10) Halgren, T. A.; Lipscomb, W. N. *Chem. Phys. Lett.* **1977**, *49*, 225.
- 581 (11) Bauer, M. S.; Strodel, B.; Fejer, S. N.; Koslover, E. F.; Wales, D. J. *J. Chem. Phys.*
582 **2010**, *132*, 054101.
- 583 (12) Wales, D. J.; Carr, J. M. *J. Chem. Theory Comput.* **2012**, *8*, 5020.

- 584 (13) Smidstrup, S.; Pedersen, A.; Stokbro, K.; Jonsson, H. *J. Chem. Phys.* **2014**, *140*,
585 214106.
- 586 (14) Sheppard, D.; Terrell, R.; Henkelman, G. *J. Chem. Phys.* **2008**, *128*, 134106.
- 587 (15) Henkelman, G.; Jonsson, H. *J. Chem. Phys.* **2000**, *113*, 9978.
- 588 (16) Henkelman, G.; Uberuaga, B.; Jonsson, H. *J. Chem. Phys.* **2000**, *113*, 9901.
- 589 (17) Weinan, E.; Ren, W. Q.; Vanden-Eijnden, E. *Phys. Rev. B* **2002**, *66*, 052301.
- 590 (18) Peters, B.; Heyden, A.; Bell, A.; Chakraborty, A. *J. Chem. Phys.* **2004**, *120*, 7877.
- 591 (19) Goodrow, A.; Bell, A. T.; Head-Gordon, M. *J. Chem. Phys.* **2008**, *129*, 174109.
- 592 (20) Behn, A.; Zimmerman, P. M.; Bell, A. T.; Head-Gordon, M. *J. Chem. Theory Comput.*
593 **2011**, *7*, 4019.
- 594 (21) Zimmerman, P. *J. Chem. Theory Comput.* **2013**, *9*, 3043.
- 595 (22) Zimmerman, P. M. *J. Chem. Phys.* **2013**, *138*, 184102.
- 596 (23) Behn, A.; Zimmerman, P. M.; Bell, A. T.; Head-Gordon, M. *J. Chem. Phys.* **2011**,
597 *135*, 224108.
- 598 (24) Chaffey-Millar, H.; Nikodem, A.; Matveev, A. V.; Krueger, S.; Roesch, N. *J. Chem.*
599 *Theory Comput.* **2012**, *8*, 777.
- 600 (25) Plessow, P. *J. Chem. Theory Comput.* **2013**, *9*, 1305.
- 601 (26) Cerjan, C. J.; Miller, W. H. *J. Chem. Phys.* **1981**, *75*, 2800.
- 602 (27) Simons, J.; Jorgensen, P.; Taylor, H.; Ozment, J. *J. Phys. Chem.* **1983**, *87*, 2745.
- 603 (28) Banerjee, A.; Adams, N.; Simons, J.; Shepard, R. *J. Phys. Chem.* **1985**, *89*, 52.
- 604 (29) Baker, J. *J. Comput. Chem.* **1986**, *7*, 385.

- 605 (30) Henkelman, G.; Jonsson, H. *J. Chem. Phys.* **1999**, *111*, 7010.
- 606 (31) Heyden, A.; Bell, A.; Keil, F. *J. Chem. Phys.* **2005**, *123*, 224101.
- 607 (32) Munro, L. J.; Wales, D. J. *Phys. Rev. B* **1999**, *59*, 3969.
- 608 (33) Machado-Charry, E.; Beland, L. K.; Caliste, D.; Genovese, L.; Deutsch, T.;
609 Mousseau, N.; Pochet, P. *J. Chem. Phys.* **2011**, *135*, 034102.
- 610 (34) Shang, C.; Liu, Z.-P. *J. Chem. Theory Comput.* **2012**, *8*, 2215.
- 611 (35) Zeng, Y.; Xiao, P.; Henkelman, G. *J. Chem. Phys.* **2014**, *140*, 044115.
- 612 (36) Peng, C. Y.; Ayala, P. Y.; Schlegel, H. B.; Frisch, M. J. *J. Comput. Chem.* **1996**, *17*,
613 49.
- 614 (37) Baker, J.; Chan, F. R. *J. Comput. Chem.* **1996**, *17*, 888.
- 615 (38) Billeter, S. R.; Turner, A. J.; Thiel, W. *Phys. Chem. Chem. Phys.* **2000**, *2*, 2177.
- 616 (39) Kato, S.; Fukui, K. *J. Am. Chem. Soc.* **1976**, *98*, 6395.
- 617 (40) Ishida, K.; Morokuma, K.; Komornicki, A. *J. Chem. Phys.* **1977**, *66*, 2153.
- 618 (41) Fukui, K. *Acc. Chem. Res.* **1981**, *14*, 363.
- 619 (42) Gonzalez, C.; Schlegel, H. B. *J. Chem. Phys.* **1989**, *90*, 2154.
- 620 (43) Taketsugu, T.; Gordon, M. S. *J. Chem. Phys.* **1995**, *103*, 10042.
- 621 (44) Wales, D.; Doye, J. *J. Phys. Chem. A* **1997**, *101*, 5111.
- 622 (45) Wales, D. J.; Scheraga, H. A. *Science* **1999**, *285*, 1368.
- 623 (46) Ohno, K.; Maeda, S. *Chem. Phys. Lett.* **2004**, *384*, 277.
- 624 (47) Maeda, S.; Ohno, K. *J. Phys. Chem. A* **2005**, *109*, 5742.

- 625 (48) Watanabe, Y.; Maeda, S.; Ohno, K. *Chem. Phys. Lett.* **2007**, *447*, 21.
- 626 (49) Ohno, K.; Maeda, S. *Phys. Scr.* **2008**, *78*, 058122.
- 627 (50) Schaefer, B.; Mohr, S.; Amsler, M.; Goedecker, S. *J. Chem. Phys.* **2014**, *140*, 214102.
- 628 (51) Zimmerman, P. M. *J. Comput. Chem.* **2013**, *34*, 1385.
- 629 (52) Rappoport, D.; Galvin, C. J. .; Zubarev, D. Y.; Aspuru-Guzik, A. *J. Chem. Theory*
630 *Computation* **2014**, *10*, 897.
- 631 (53) Maeda, S.; Morokuma, K. *J. Chem. Theory Comput.* **2011**, *7*, 2335.
- 632 (54) Maeda, S.; Ohno, K.; Morokuma, K. *Phys. Chem. Chem. Phys.* **2013**, *15*, 3683.
- 633 (55) Laio, A.; Parrinello, M. *Proc. Natl. Acad. Sci. USA* **2002**, *99*, 12562.
- 634 (56) Carter, E. A.; Ciccotti, G.; Hynes, J. T.; Kapral, R. *Chem. Phys. Lett.* **1989**, *156*, 472.
- 635 (57) Voter, A. *Phys. Rev. Lett.* **1997**, *78*, 3908.
- 636 (58) Iannuzzi, M.; Laio, A.; Parrinello, M. *Phys. Rev. Lett.* **2003**, *90*, 238302.
- 637 (59) Ensing, B.; De Vivo, M.; Liu, Z.; Moore, P.; Klein, M. *Acc. Chem. Res.* **2006**, *39*, 73.
- 638 (60) Pietrucci, F.; Andreoni, W. *Phys. Rev. Lett.* **2011**, *107*, 085504.
- 639 (61) Sorensen, M.; Voter, A. *J. Chem. Phys.* **2000**, *112*, 9599.
- 640 (62) Martinez-Nunez, E. *J. Comput. Chem.* **2015**, *36*, 222.
- 641 (63) Bernasconi, M.; Chiarotti, G. L.; Focher, P.; Parrinello, M.; Tosatti, E. *Phys. Rev. Lett.*
642 **1997**, *78*, 2008.
- 643 (64) Goldman, N.; Reed, E. J.; Fried, L. E.; Kuo, I.-F. W.; Maiti, A. *Nature Chem.* **2010**,
644 *2*, 949.

- 645 (65) Wang, L.-P.; Titov, A.; McGibbon, R.; Liu, F.; Pande, V. S.; Martinez, T. J. *Nature*
646 *Chem.* **2014**, *6*, 1044.
- 647 (66) Titov, A. V.; Ufimtsev, I. S.; Luehr, N.; Martinez, T. J. *J. Chem. Theory Comput.*
648 **2013**, *9*, 213.
- 649 (67) Ufimtsev, I. S.; Martinez, T. J. *J. Chem. Theory Comput.* **2009**, *5*, 2619.
- 650 (68) van Duin, A. C. T.; Dasgupta, S.; Lorant, F.; Goddard III, W. A. *J. Phys. Chem. A*
651 **2001**, *105*, 9396.
- 652 (69) Elstner, M.; Porezag, D.; Jungnickel, G.; Elsner, J.; Haugk, M.; Frauenheim, T.;
653 Suhai, S.; Seifert, G. *Phys. Rev. B* **1998**, *58*, 7260.
- 654 (70) Cordero, B.; Gomez, V.; Platero-Prats, A. E.; Reves, M.; Echeverria, J.; Cremades, E.;
655 Barragan, F.; Alvarez, S. *Dalton Trans.* **2008**, *37*, 2832.
- 656 (71) Russell, S.; Norvig, P. *Artificial Intelligence: A Modern Approach*, 3rd ed.; Pearson
657 Education: Upper Saddle River, New Jersey 07458, 2010.
- 658 (72) Shao, Y.; Molnar, L. F.; Jung, Y.; Kussmann, J.; Ochsenfeld, C.; Brown, S. T.;
659 Gilbert, A. T. B.; Slipchenko, L. V.; Levchenko, S. V.; O'Neill, D. P.; Jr., R. A. D.;
660 Lochan, R. C.; Wang, T.; Beran, G. J. O.; Besley, N. A.; Herbert, J. M.; Lin, C. Y.;
661 Voorhis, T. V.; Chien, S. H.; Sodt, A.; Steele, R. P.; Rassolov, V. A.; Maslen, P. E.;
662 Korambath, P. P.; Adamson, R. D.; Austin, B.; Baker, J.; Byrd, E. F. C.; Dachsel, H.;
663 Doerksen, R.; Dreuw, A.; Dunietz, B. D.; Dutoi, A. D.; Furlani, T. R.; Gwaltney, S. R.;
664 Heyden, A.; Hirata, S.; Hsu, C.-P.; Kedziora, G.; Khalliulin, R. Z.; Klunzinger, P.;
665 Lee, A. M.; Lee, M. S.; Liang, W.; Lotan, I.; Nair, N.; Peters, B.; Proynov, E. I.; Pieni-
666 azek, P. A.; Rhee, Y. M.; Ritchie, J.; Rosta, E.; Sherrill, C. D.; Simmonett, A. C.; Sub-
667 otnik, J. E.; III, H. L. W.; Zhang, W.; Bell, A. T.; Chakraborty, A. K.; Chipman, D. M.;
668 Keil, F. J.; Warshel, A.; Hehre, W. J.; III, H. F. S.; Kong, J.; Krylov, A. I.; Gill, P.
669 M. W.; Head-Gordon, M. *Phys. Chem. Chem. Phys.* **2006**, *8*, 3172.

- 670 (73) Suits, A. G. *Acc. Chem. Res.* **2008**, *41*, 873.
- 671 (74) Bowman, J. M. *Mol. Phys.* **2014**, *112*, 2516.
- 672 (75) Roos, B. O. *Adv. Chem. Phys.* **1987**, *69*, 399.
- 673 (76) Hohenstein, E. G.; Luehr, N.; Ufimtsev, I. S.; Martinez, T. J. *J. Chem. Phys.* **2015**,
674 *142*, 224103.
- 675 (77) Hohenstein, E. G.; Bouduban, M. E. F.; Song, C.; Luehr, N.; Ufimtsev, I. S.; Mar-
676 tinez, T. J. *J. Chem. Phys.* **2015**, *143*, 014111.
- 677 (78) Shu, Y.; Hohenstein, E. G.; Levine, B. G. *J. Chem. Phys.* **2015**, *142*, 024102.
- 678 (79) Tomasi, J.; Mennucci, B.; Cammi, R. *Chem. Rev.* **2005**, *105*, 2999.
- 679 (80) Liu, F.; Luehr, N.; Kulik, H. J.; Martinez, T. J. *J. Chem. Theory Comput.* **2015**, *11*,
680 3131.
- 681 (81) Lange, A. W.; Herbert, J. M. *J. Chem. Phys.* **2010**, *133*, 244111.

682 **Graphical TOC Entry**

683

