

From high-dimensional committors to reactive insights

Nils E. Strand,^{1,*} Schuyler B. Nicholson,¹ Hadrien Vroylandt,^{1,†} and Todd R. Gingrich^{1,‡}

¹*Department of Chemistry, Northwestern University,
2145 Sheridan Road, Evanston, Illinois 60208, USA*

Transition path theory (TPT) offers a powerful formalism for extracting the rate and mechanism of rare dynamical transitions between metastable states. Most applications of TPT either focus on systems with modestly sized state spaces or use collective variables to try to tame the curse of dimensionality. Increasingly, expressive function approximators like neural networks and tensor networks have shown promise in computing the central object of TPT, the committor function, even in very high-dimensional systems. That progress prompts our consideration of how one could use such a high-dimensional function to extract mechanistic insight. Here, we present and illustrate a straightforward but powerful way to track how individual dynamical coordinates evolve during a reactive event. The strategy, which involves marginalizing the reactive ensemble, naturally captures the evolution of the dynamical coordinate's distribution, not just its mean reactive behavior.

I. INTRODUCTION

Biophysical systems frequently involve dynamics that is both high-dimensional and stochastic [1–7]. When those dynamical processes relax into an equilibrium, it is possible to study the stable states in terms of thermodynamics without reference to the dynamics. However, many biophysical processes operate away from equilibrium, and in that regime, it is especially crucial to understand the kinetics. Studying Markovian models is a well-developed route to analyzing those high-dimensional stochastic kinetics, a route common to the chemical master equation [8], Langevin dynamics [9], and Markov State Models (MSMs) [10, 11]. It is often the case that these Markovian models exhibit slow-time-scale transitions between different regions of a configuration space, transitions one would associate with barrier crossing in an equilibrium setting. Owing to the chemical physics history, we generically refer to changes from one set of metastable configurations \mathcal{A} to another \mathcal{B} as a reaction or a reactive event. Examples of such reactions are transitions between metastable states in gene regulatory networks and chemical reaction networks [12–14]. Because the configuration space can become astronomically vast, one often seeks a coarse-grained description of the kinetics: What are the long-lived metastable regions of the configuration space, what are the timescales for reactions, and what is the mechanism of the reaction? That mechanism is particularly desirable, as it is easier to design ways to modify the reaction rate if one knows *how* the reaction typically proceeds.

The most straightforward approach to learning the mechanism involves generating and watching ensembles of representative reactive trajectories [15, 16] to form impressions of how those representative trajectories

progress from \mathcal{A} to \mathcal{B} . Due to a separation of timescales between the typical residence time in metastable states and the transition time [17], it can be impractical to directly simulate and watch the large number of required trajectories. Enhanced sampling methods such as transition path sampling (TPS) [18–21] and forward flux sampling (FFS) [6, 22, 23] can offer more efficient ways to generate the ensemble of reactive trajectories, but even when the ensemble can be sampled, the results are still high-dimensional, making them nontrivial to interpret. What is needed is a low-dimensional representation of the resulting mechanism from the high-dimensional reactive trajectories.

To avoid these challenges, one can parameterize the progress not by time but by a one-dimensional reaction coordinate. The reaction coordinate can be thought of as a many-to-one mapping from microstate \mathbf{x} onto a single variable measuring progress along the reaction, and the best choice for such a progress coordinate is known to be the so-called committor function, $q(\mathbf{x})$ [24, 25]. Transition path theory (TPT) [24, 26–28] provides explicit expressions to compute $q(\mathbf{x})$ in terms of a generator of dynamics, but the cost of directly performing such computations rapidly increases with the number of microstates. Due to the curse of dimensionality, it is common for problems of interest to have astronomically many microstates. In these cases, the most common way that committor functions have been used for complex systems is to avoid $q(\mathbf{x})$ and instead compute a committor function defined over a low-dimensional (often one-dimensional) collective variable \mathbf{y} , which is a function of \mathbf{x} . This committor $q(\mathbf{y})$ is practically computed by sampling. For example, many trajectories can be initialized with a particular value of \mathbf{y} and then propagated until they reach either \mathcal{A} or \mathcal{B} , with $q(\mathbf{y})$ being the probability of first reaching \mathcal{B} . Approaches built around committors of one or more collective variables have been productive [29–40], but the approaches typically require the choice of good collective variables upfront. A significant body of research has developed strategies to identify and optimize those “good” collective variables, ideally finding a \mathbf{y} that resembles the committor itself [19, 41–50]. More recently, basis ex-

* Now at James Franck Institute, University of Chicago, Chicago, IL 60637, USA

† Now at CERMICS, École des Ponts ParisTech.

‡ todd.gingrich@northwestern.edu.

pansions [51–54], neural networks [32, 55–61] and tensor networks [62] have been used to estimate $q(\mathbf{y})$ from sampled trajectories even when \mathbf{y} is quite high dimensional. Those advances pair nicely with strategies to extract the mechanism of reactive events in the collective variable space from $q(\mathbf{y})$ [63–65]. In particular, Ref. [54] has used that $q(\mathbf{y})$ to inspect how the steady-state distribution of certain collective variables varies as a function of reaction progress.

Suppose, by contrast, that one could discard the collective variables altogether and it were practical to solve for the full-dimensional committor $q(\mathbf{x})$. Here, we introduce and illustrate that for discrete systems such as well-mixed chemical reaction networks, it is indeed numerically practical to compute $q(\mathbf{x})$ and extract mechanistic insight. For these problems, the strategy does not even require trajectory sampling. The key idea of this paper is that access to the full-dimensional committor $q(\mathbf{x})$ allows one to inspect how each dynamical coordinate $x_i \in \mathbf{x}$ evolves as a function of reaction progress. Crucially, this approach retains a distribution over x_i , not just the mean behavior, allowing the approach to directly reveal the presence of multiple reactive pathways.

The methodology is built upon the reactive ensemble of TPT, which gives the density $\rho^{AB}(\mathbf{x})$ of occupying a microstate \mathbf{x} given that the system is in the midst of transitioning from \mathcal{A} to \mathcal{B} . For each degree of freedom x_i , $1 \leq i \leq D$, we compute a two-dimensional distribution formed from the reactive ensemble by marginalizing over all other degrees of freedom:

$$\begin{aligned} \rho^{AB}(x_i, q) \\ = \int dx_1 \cdots dx_{i-1} dx_{i+1} \cdots dx_D \rho^{AB}(\mathbf{x}) \delta(q(\mathbf{x}) - q), \end{aligned} \quad (1)$$

with the δ function serving to pick out how far the reaction had progressed toward \mathcal{B} . The marginal $\rho^{AB}(x_i, q)$ highlights the single coordinate x_i , but it retains the influence of the other coordinates only in so far as they impact the progress coordinate q . In this way, one can view how the distribution for each x_i evolves during the reaction process, parameterized by q . The approach has the flavor of the so-called violin plots of Ref. [54], but it computes the q dependence of the reactive, rather than steady-state, ensemble. The reactive ensemble might be a preferred choice if the steady state has overwhelming probability within the regions \mathcal{A} and \mathcal{B} , making the steady-state probabilities between the two metastable basins difficult to visualize.

We illustrate the idea with two example problems, both of which admit a direct computation of q over a discrete state space. First, we demonstrate the approach for a discretized two-dimensional (2D) diffusion problem where the explicit calculation of the committor has previously been studied [28, 66, 67]. Though only two-dimensional, this problem illustrates the approach and emphasizes that it can naturally highlight when multiple distinct pathways meaningfully contribute to reac-

tive events. Second, we move to a situation with too many degrees of freedom to straightforwardly plot $q(\mathbf{x})$, a gene toggle switch (GTS) model [6] with two metastable states emerging from stochastic chemical kinetics of seven chemical species. Calculating the committor for the GTS model is more complicated than most literature toy problems since we consider a GTS model with several million microstates, many more than the coarse-grained models typically used for transition path analysis [63, 68–73]. Using sparse linear algebra methods, we compute q and show how it can be used to extract the reaction mechanism one species at a time.

II. TRANSITION PATH THEORY

A. Standard formulation

Our work builds upon TPT, so we start by reviewing its main results for discrete-state continuous-time Markov dynamics [24]. One can choose a canonical ordering of microstates so a many-body microstate \mathbf{x} is labeled by the single index i . Let W_{ij} denote the rate or probability per unit time of transitions from the j^{th} into the i^{th} microstate. Conservation of probability is enforced because the diagonal elements W_{ii} are chosen such that $\sum_i W_{ij} = 0$. Without loss of generality, we assume that it is possible to reach any microstate from any other microstate in a finite number of transitions; that is, W is irreducible. In the long-time limit, the microstate i is visited with steady-state probability π_i . That distribution follows simply from the matrix W as the solution to

$$W\boldsymbol{\pi} = 0, \quad (2)$$

where $\boldsymbol{\pi}$ is the vector of steady-state probabilities for each microstate.

TPT partitions the space of microstates \mathcal{S} into three regions: \mathcal{A} , \mathcal{B} , and $(\mathcal{A} \cup \mathcal{B})^c$, where the superscript c is the complement of the set. The aim is to describe properties of the Markov dynamics within $(\mathcal{A} \cup \mathcal{B})^c$ conditioned upon starting in \mathcal{A} and ending in \mathcal{B} , without having first returned to \mathcal{A} . This conditioned process is of special physical interest when \mathcal{A} and \mathcal{B} are metastable states and trajectories pass through $(\mathcal{A} \cup \mathcal{B})^c$ rarely. The rare transitions are then viewed as reactions. A motivating goal for TPT was to compute the reaction rate k_{AB} from the Markov rate operator W . It has been shown that this reaction rate is expressed compactly in terms of the committor function, specifically the forward committor function, which we distinguish with a superscript $+$. For the discrete state space, we define the vector \mathbf{q}^+ whose element q_i^+ is the probability that a trajectory initiated in state i will reach \mathcal{B} before \mathcal{A} . The forward committor solves the Dirichlet boundary value problem

$$\begin{cases} \sum_{j \in \mathcal{S}} q_j^+ W_{ji} = 0, & \forall i \in (\mathcal{A} \cup \mathcal{B})^c \\ q_i^+ = 0, & \forall i \in \mathcal{A} \\ q_i^+ = 1, & \forall i \in \mathcal{B} \end{cases}. \quad (3)$$

Practically, it is convenient to cast that problem as the linear equation

$$U\mathbf{q}^+ = \mathbf{v}, \quad (4)$$

where U is a square matrix with elements $U_{ji} = W_{ij}$, for all $i, j \notin \mathcal{A} \cup \mathcal{B}$, and \mathbf{v} is a vector with elements $v_j = -\sum_{i \in \mathcal{B}} W_{ij}$, for all $j \in (\mathcal{A} \cup \mathcal{B})^c$, and zero otherwise. In words, U is the transpose of the submatrix of W corresponding to the reactive region, that is, the sites not in \mathcal{A} or \mathcal{B} . Multiplying the forward committor vector \mathbf{q}^+ from the left by U results in \mathbf{v} , a vector whose element $i \notin \mathcal{A} \cup \mathcal{B}$ is minus the sum of the rates leaving i and entering \mathcal{B} .

TPT defines the backward committor \mathbf{q}^- in a manner analogous to \mathbf{q}^+ ; q_i^- is the probability of being in i given that the system last occupied \mathcal{A} before \mathcal{B} . The backward committor relies on the time-reversed process, characterized by a rate matrix \tilde{W} . The off-diagonal elements of this matrix are given by $\tilde{W}_{ij} = W_{ji}\pi_i/\pi_j$, and the diagonal elements are $\tilde{W}_{ii} = -\sum_{j \neq i} \tilde{W}_{ji}$ [28]. The boundary value problem for the backward committor,

$$\begin{cases} \sum_{j \in \mathcal{S}} q_j^- \tilde{W}_{ji} = 0, & \forall i \in (\mathcal{A} \cup \mathcal{B})^c \\ q_i^- = 1, & \forall \mathbf{x} \in \mathcal{A} \\ q_i^- = 0, & \forall \mathbf{x} \in \mathcal{B} \end{cases}, \quad (5)$$

leads to the linear equation

$$\tilde{U}\mathbf{q}^- = \tilde{\mathbf{v}}, \quad (6)$$

where $\tilde{U}_{ij} = \tilde{W}_{ji}$, for all $i, j \notin \mathcal{A} \cup \mathcal{B}$ and $\tilde{v}_i = -\sum_{k \in \mathcal{A}} \tilde{W}_{ki}$, for all $i \in (\mathcal{A} \cup \mathcal{B})^c$. For reversible systems, $W = \tilde{W}$, and forward and backward committors are trivially related as $q_i^- = 1 - q_i^+$. The calculations in this paper involve Markov dynamics with a W that breaks detailed balance, requiring Eqs. (4) and (6) to be solved independently. By additionally solving for the steady-state distribution $\boldsymbol{\pi}$, one can then construct the reactive probability

$$P_i^{AB} = \pi_i q_i^+ q_i^-. \quad (7)$$

This P_i^{AB} is the probability that a reactive trajectory occupies discrete microstate i [28]. If these microstates come from a discretization of a continuous problem, then the reactive density is $\rho^{AB}(\mathbf{x}) = P_i^{AB}/V$, where i is the index for microstate \mathbf{x} and V is the volume element of each discretized cell.

B. Re-expression for large-scale systems breaking detailed balance

When detailed balance is not satisfied, as in our second example problem, \mathbf{q}^- does not follow directly from \mathbf{q}^+ . In those cases, \mathbf{q}^- could in principle be found by solving Eq. (6), but that approach is impractical for large systems. Elements of the time-reversed rate matrix can

suffer from numerical instability due to states with vanishingly small steady-state probabilities entering into the denominator of $\tilde{W}_{ij} = W_{ji}\pi_i/\pi_j$. Note, however, that we can directly solve for the vector \mathbf{r} whose elements $r_i = \pi_i q_i^-$ appear in the expression for the reactive density, Eq. (7). In Appendix B, we demonstrate that solving Eq. 6 is equivalent to solving the linear equation

$$U^T \mathbf{r} = \mathbf{s}, \quad (8)$$

where $s_i = -\sum_{k \in \mathcal{A}} W_{ik}\pi_k$ for all $i \in (\mathcal{A} \cup \mathcal{B})^c$. Notice that unlike Eq. (6), Eq. (8) requires a transpose of the operator U used to solve the forward committor, not the numerically problematic time-reversal. The transformation does not give a free lunch in that the vector \mathbf{s} cannot be constructed as simply as \mathbf{v} and $\tilde{\mathbf{v}}$. Instead, \mathbf{s} requires knowledge of the steady state $\boldsymbol{\pi}$. We therefore solve for \mathbf{r} in two stages. First, we find $\boldsymbol{\pi}$ by applying Arnoldi iteration to the eigenvalue problem in Eq. (2). With that $\boldsymbol{\pi}$, we construct \mathbf{s} and use general minimum residual (GMRES) iterations to solve Eqs. (4) and (8) for \mathbf{q}^+ and \mathbf{r} . Combining the two, we obtain the reactive distribution as $P_i^{AB} = q_i^+ r_i$.

III. RESULTS

A. Two-dimensional diffusion on a metastable landscape

Before breaking detailed balance or considering high-dimensional systems, it is useful to discuss a simpler prototypical minimal example, that of two-dimensional diffusion on a metastable landscape [28]. For this example, the microstates \mathbf{x} are defined by two coordinates, x and y . The system evolves on the three-well energy landscape

$$\begin{aligned} V(x, y) = & 3e^{-x^2 - (y - \frac{1}{3})^2} - 3e^{-x^2 - (y - \frac{5}{3})^2} \\ & - 5e^{-(x-1)^2 - y^2} - 5e^{-(x+1)^2 - y^2} \\ & + \frac{x^4}{5} + \frac{(y - \frac{1}{3})^4}{5}, \end{aligned} \quad (9)$$

according to an overdamped Langevin dynamics with a gradient force and a random force $\boldsymbol{\xi}$ of thermal origin:

$$\dot{\mathbf{x}} = -\nabla V(\mathbf{x}) + \boldsymbol{\xi}. \quad (10)$$

With inverse temperature β , the white noise satisfies $\langle \xi_i(t) \xi_j(t') \rangle = 2\beta^{-1} \delta_{ij} \delta(t - t')$. The energy landscape was constructed to have metastable basins (see Fig. 1a), and the standard problem is to describe the rare dynamical path that causes the system to transition from one of those basins to the other. This particular toy problem is a useful starting point because the state space can be discretized onto a 200×200 grid such that the corresponding linear equation, Eq. (4), for the forward committor can be solved [28]. The fineness of the grid one uses depends on the relative size of the deterministic

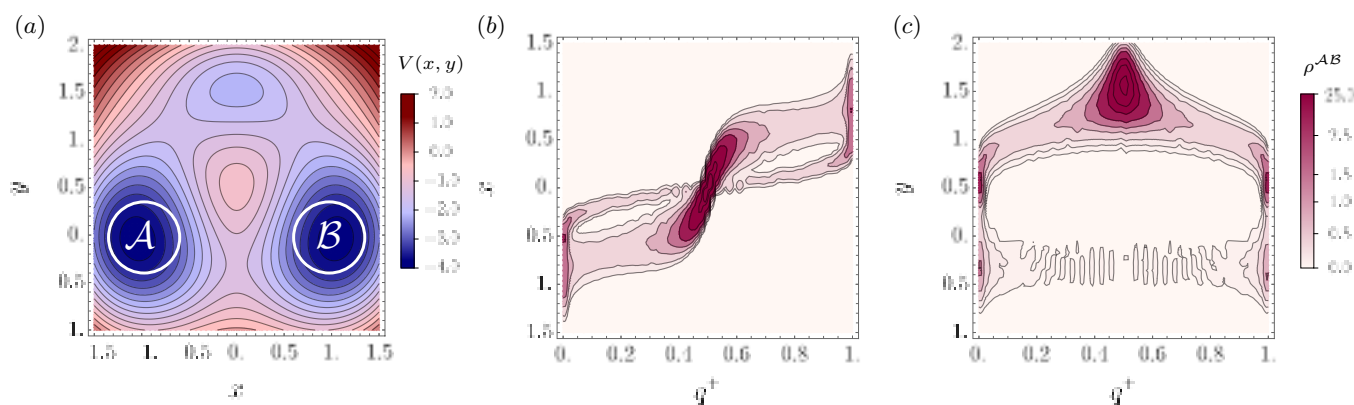


FIG. 1. Detail-balanced dynamics of a single thermal particles. Following Ref. [28], we study overdamped dynamics on an energy landscape with two deep wells and seek information about the mechanism passing from reactants \mathcal{A} to products \mathcal{B} . Upon discretizing the continuous state space on a grid, sparse linear algebra methods give π , \mathbf{q}^+ , and $\mathbf{q}^- = 1 - \mathbf{q}^+$ for that discretized problem. The reactive density $\rho^{AB}(\mathbf{x})$ can be marginalized as in Eq. (1) to reveal the distribution of each coordinate (x and y) as a function of reaction progress q^+ . Plots of $\rho^{AB}(x, q^+)$ and $\rho^{AB}(y, q^+)$ reflect the distribution of outcomes for coordinate x and y , respectively, where the state of the other coordinate is considered only to the extent that it impacts the reaction progress q^+ . Retaining information about the statistical ensemble naturally reveals the presence of multiple reactive pathways. For cleaner visualization, the reactive ensemble densities were smoothed with Gaussian kernels and bandwidths obeying Silverman's rule [74] and colored with a nonlinear hyperbolic tangent colorbar that enhances the resolution of the low-probability regions.

forces and the stochastic white-noise forces that generate diffusion. At a minimum, a reasonable discretization must be sufficiently fine that the rate operator has non-negative off-diagonal elements [75]. With \mathbf{x} being only two-dimensional, the solution $q^+(\mathbf{x})$ can be plotted to offer a clear visual for how the reactions proceed from one basin to the other. The landscape has two distinct pathways along which transitions can occur, and a plot of $q^+(\mathbf{x})$ shows which pathway dominates [28].

The challenge we set out to address is how one can use $q^+(\mathbf{x})$ to describe the typical reaction mechanism when \mathbf{x} is too high-dimensional to plot. Because of the curse of dimensionality, we need to consider the components of \mathbf{x} one-by-one, inspecting how each one advances, and yet the two-dimensional diffusion problem highlights the difficulty of decoupling those degrees of freedom. A plot of $q^+(x, y)$ contains information about how the correlated motion of x and y can conspire to advance the reaction. If we simply neglected the information about y , we would lose those correlations. Figure 1 shows how our proposal, Eq. (1), can capture how each coordinate evolves one-by-one, even resolving the multiple reaction pathways. To highlight this capability, we solved for $q^+(\mathbf{x})$ at inverse temperature $\beta = 4$, a value at which both pathways contribute meaningfully to transitions. As a function of reaction progress (q^+), we monitor how the distribution of each coordinate evolves, revealing distinctly bimodal distributions that form two channels in Fig. 1. The plot in Fig. 1c, for example, reflects that progress can emerge either when y increases to a shallow intermediate basins around $(x, y) = (0, 1.5)$ or by holding $y \approx 0$ and letting x do all of the work by climbing up and over a single saddle. Because the plots show a *distribution*, not just

a mean, at each value of q^+ , they contain rich information about how important the coordinate's motion is for enabling reaction progress. The channels stretch along the vertical direction when reaction progress is relatively insensitive to the precise value of the coordinate, and they narrow when the coordinate is strongly driving the reaction. For a two-dimensional \mathbf{x} we do not mean to suggest that plotting $\rho^{AB}(x, q^+)$ and $\rho^{AB}(y, q^+)$ is simpler to parse than a plot of $q^+(x, y)$. Rather, our point is that rich insights about the reaction mechanism can be extracted by these collections of plots, which remain computable and interpretable even when the dimensionality grows. We emphasize this point with the second example.

B. Bistable switching in a seven-species gene toggle switch

1. The GTS Model

A paradigmatic example of bistable transitions in higher dimensions is provided by the chemical master equation (CME) for the stochastic dynamics of a gene toggle switch (GTS) [6]. The GTS model that we studied was constructed to describe the fluctuating copy numbers of two proteins, A and B. A single piece of DNA, denoted in the model by O, containing genes for A and B provides routes to increase the copy numbers through protein synthesis, but the copy numbers can also decrease via protein degradation. The two genes mutually suppress each other, e.g., increasing the number of A decreases the pro-

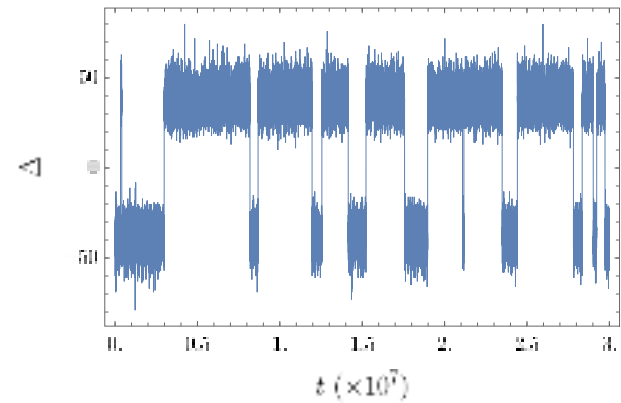
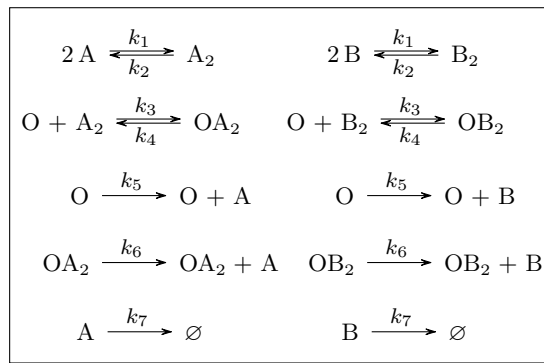


FIG. 2. **Gene toggle switch (GTS) model.** Left: Set of chemical reactions that make up the GTS model [6], a chemical reaction network model with 7 species that perform 14 stochastic reactions with 7 distinct reaction rates, k_1 through k_7 . The system involves a single piece of DNA (species O) that can synthesize protein A or protein B. Those proteins can dimerize and the dimers can bind to DNA as a promoter that suppresses the production of the other protein. The result is a bistable switch that toggles between an A-rich and B-rich state. Right: Which state the system occupies is well captured by the order parameter Δ that counts how many more A proteins there are than B proteins (in monomer, dimer, or bound forms). A brute-force stochastic simulation gives a Monte Carlo realization of a trajectory, illustrating the stochastic switching observed for the numerical values of the reaction rates reported in the main text. Fig. 3 views this same process from the perspective of a statistical ensemble by employing transition path theory (TPT).

duction rate of B. Consequently, typical microstates involve either a high number of A or a high number of B, with rare stochastic fluctuations toggling between the metastable states.

The specific GTS model we study involves seven chemical species and fourteen reactions (see Fig. 2). The model allows for reversible dimerization of A and B to produce A_2 and B_2 . Each dimer can also reversibly bind to the DNA to give OA_2 and OB_2 . The bound dimer acts as a promoter, so OA_2 prompts the synthesis of more copies of protein A without similarly prompting the synthesis of B. In the absence of a bound promoter, O is equally likely to synthesize A and B. Finally, both proteins have an irreversible degradation process. Figure 2 labels the rates for each of the fourteen elementary reactions by k_1 through k_7 , assuming a symmetry between the kinetics of A and B. We follow Ref. [6], setting $k_1 = k_2 = k_3 = 5$, $k_4 = k_5 = k_6 = 1$, and $k_7 = 0.25$. The symmetry between A and B is spontaneously broken by the fluctuating dynamics, and the imbalance is monitored by the order parameter $\Delta \equiv n_A + 2n_{A_2} + 2n_{OA_2} - n_B - 2n_{B_2} - 2n_{OB_2}$, where n_γ denotes the number of species γ . A microstate for the GTS model is then given by $(n_A, n_{A_2}, n_{OA_2}, n_O, n_{OB_2}, n_{B_2}, n_B)$. The representative stochastic trajectory of Fig. 2 shows that one can define a metastable \mathcal{A} region by $\Delta \geq 25$ and a metastable \mathcal{B} region by $\Delta \leq -25$. Although the vast majority of the time is spent within either basin, we are primarily interested in the behavior of trajectories leaving \mathcal{A} and entering \mathcal{B} .

Like in the first example, the stochastic dynamics of the GTS is described by a Markovian jump process from one microstate to another. The two-dimensional diffu-

sion required discretization onto a grid, but the states of the GTS naturally occupy a seven-dimensional lattice, one dimension per chemical species (A, A_2 , OA_2 , O, OB_2 , B_2 , and B). None of the reactions can make or destroy the DNA, so there is a constant of motion: $n_{OA_2} + n_O + n_{OB_2} = 1$. That constraint restricts species OA_2 , O, and OB_2 to each be present with zero or one copy. In contrast, the copy number of the A and B monomers and dimers can in principle increase without bound. In practice, the degradation rate k_7 ensures that there is some maximum copy number, M , above which the dynamics is exceedingly unlikely to sample. Appendix A shows that the truncation at $M = 30$ did not appreciably influence the reactive trajectories. This choice of M is somewhat larger than what one might intuitively deduce from the positions of the distribution peaks, as shown in Fig. 4, because the shape of the distribution is non-trivially influenced by rare configurations in the tails of the distribution. Since we restrict n_A, n_{A_2}, n_B , and n_{B_2} to each be between 0 and M , the state must be one of $3(M+1)^4$ microstates which comes to nearly 2.8 million for $M = 30$. Although it is not entirely trivial to converge such a large vector, it is possible because CMEs naturally support a sparse representation of the Markov operator W .

Provided that the number of microstates is sufficiently modest that they can be practically enumerated, constructing the sparse matrix for W is straightforward. For the i^{th} microstate, one loops over the reactions in Fig. 2 to identify the microstate index j that would result if that reaction were to fire. To the sparse matrix W , one adds an element $W_{ji} = \alpha$, where the so-called propensity α is the rate of reaction k_r times a combinatorial factor count-

ing how many distinct copies of the species could have executed the reaction. For example, if the microstate i had 5 A monomers and 7 A₂ dimers, then the reaction $2A \xrightarrow{k_1} A_2$ would have a rate $20k_1$ of mapping to the microstate with 3 A monomers and 8 A₂ dimers. Here, we have 20 distinct ways that two of the five A's could have participated in the reaction. Any reaction that would have increased the copy number to exceed M would be ignored. Once all nonzero off-diagonal elements of W are identified, the diagonal elements are set to $W_{ii} = -\sum_{j \neq i} W_{ij}$ to enforce conservation of probability.

The above procedure is conceptually simple and happens to be computationally tractable for this system size, but there is a more elegant way to build W that also extends to CMEs with astronomically large numbers of microstates. The alternative approach leverages the Doi-Peliti (DP) formalism to represent W in terms of raising and lowering operators that act on each chemical species [76–78]. For the GTS problem, one can arrive at all of our results without the DP formalism, but we envision extensions to larger state spaces such that vectors like \mathbf{q}^+ cannot be explicitly computed, but are rather approximated by a tensor network. In those very large state spaces, looping over microstates is not possible and W must be built using the DP formalism. Anticipating this necessity, we describe the DP representation of the GTS model in Appendix C.

Having built W , Eq. (2), Eq. (4), and Eq. (8) are solved by routine sparse linear algebra methods. For $M = 30$, the convergence of $\boldsymbol{\pi}$ via about 5×10^4 Arnoldi iterations is the most expensive step, requiring approximately ten hours of serial runtime on a single processor. We needed 10^6 GMRES iterations to converge \mathbf{q}^+ and \mathbf{r} in 3 hours and 2 hours, respectively. Because the model has been a playground for advanced sampling algorithms, it is tempting to compare the timings of the rate calculations. In Appendix A, we show agreement between the $M = 30$ TPT rate calculation and a simple sampling rate estimate that required approximately 1500 serial CPU hours. FFS is known to accelerate similar calculations by a factor of 40-90 [6], reflecting that the TPT calculation costs the same order of magnitude as an FFS rate estimate. This comparison suggests that the TPT calculation is in the same ballpark as enhanced sampling algorithms, although we highlight that the TPT calculation yields not only the rate. At the same expense, the TPT approach also gives the committor, which we now use to analyze the mechanism.

2. Analyzing the committor

Since the GTS model breaks detailed balance, we computed $q^+(\mathbf{x})$, and $r(\mathbf{x})$ to obtain the reactive density P^{AB} of Eq. (7). Fig. 3 shows how the distribution for the number of each species evolves as a function of reaction progress q^+ . Those distributions, $P^{AB}(n_A, q^+)$,

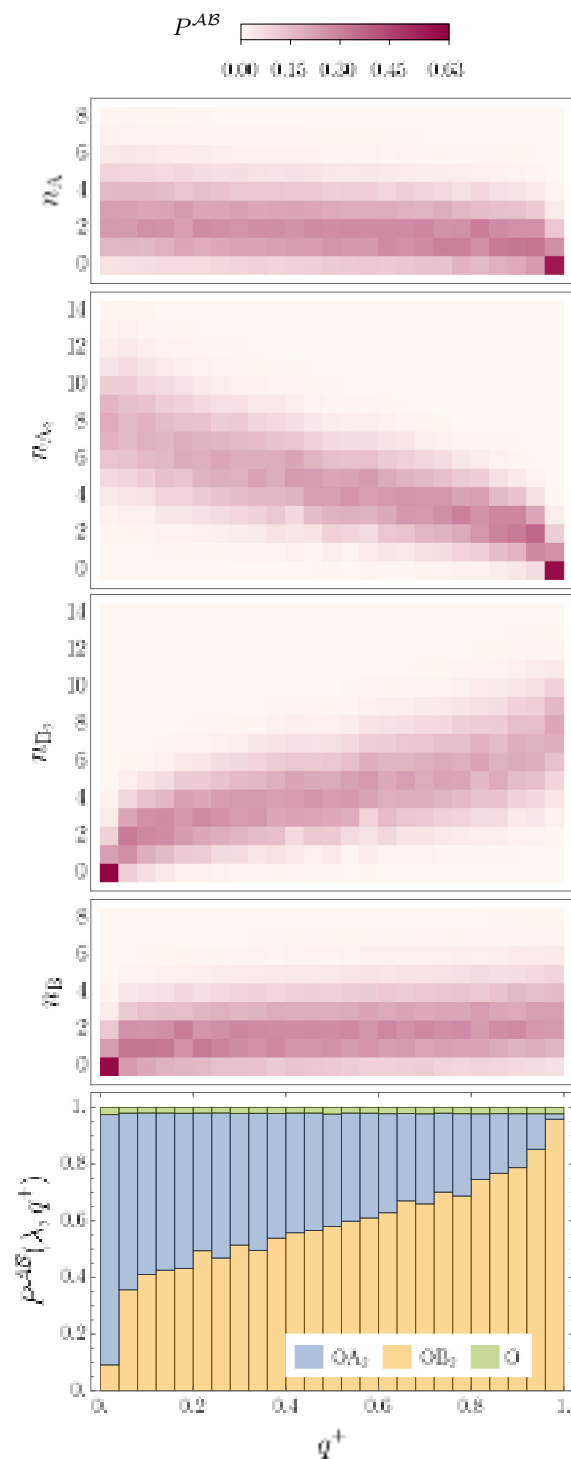


FIG. 3. GTS reactive ensemble. For the GTS model, the reaction coordinate q^+ is explicitly computed with sparse linear algebra methods, allowing the distribution for each chemical species to be tracked as a function of reaction progress. Some species undergo significant changes in their average count, while others have more subtle changes in higher moments of their distribution. Counts of those chemical species are naturally discretized, but q^+ is reported with bins of width 0.04 to aid in visualization. As marginals for the *reactive* ensemble, these plots highlight typical reactive behavior moving from the edge of region \mathcal{A} to the edge of \mathcal{B} . Configurations at those boundaries are not representative of typical steady-state behavior, which is dominated by \mathcal{A} and \mathcal{B} .

$P^{AB}(n_{A_2}, q^+)$, $P^{AB}(n_B, q^+)$, and $P^{AB}(n_{B_2}, q^+)$, follow from the marginalizations of Eq. (1), computed by discretizing q^+ with bins of width 0.04. The same procedure could also produce the reactive density for three different DNA states n_O , n_{OA_2} , and n_{OB_2} . Since DNA must be in one and only one of these three states, it is more revealing to construct a new variable λ that records which of the three states the DNA is in. That distribution over λ states then follows from a corresponding $P^{AB}(\lambda, q^+)$.

The five plots in Fig. 3 collectively tell the story of how the elementary reactions of Fig. 2 collude together to allow the system to transit from A-rich to B-rich microstates. Perhaps the clearest feature of the plots is the fact that the probability of finding the DNA in the O state is very small and completely insensitive to q^+ . The calculations therefore show that the DNA will typically be bound to a dimer and the reaction proceeds by switching from a bound A_2 to a bound B_2 . However, the reaction is not “halfway done” once the DNA flips from OA_2 to OB_2 . The bottom plot shows that $q^+ \approx 0.35$ when OA_2 and OB_2 are equally likely in the reactive ensemble. To push q^+ beyond 0.5, it is also important that a sufficiently large population of B_2 is built up, serving as a memory that prevents a rapid backsliding into the OA_2 state. Plots of $P^{AB}(n_B, q^+)$ and $P^{AB}(n_{B_2}, q^+)$ show that monomer and dimer play distinct roles. DNA in the OB_2 state produces only B, allowing for a buildup of monomer, but $P^{AB}(n_B, q^+)$ shows that the monomer does not appreciably build up over the course of the reaction. While the distribution over n_B subtly shifts as a function of q^+ , it is always rare to see much more than 4 B molecules. The relatively uniform fluctuations in n_B reflect that the number of B molecules is a poor proxy for the progress along the reaction coordinate. $P^{AB}(n_{B_2}, q^+)$ shows that it is instead the population of the dimer B_2 that drives the progress to make the OB_2 toward a stable B-rich state, a conclusion that follows from the drift in the peak of the n_{B_2} distribution as q^+ grows.

IV. DISCUSSION

In this work, we have outlined and illustrated an approach to capture the mechanism of transitions between two regions of very high-dimensional complex systems. Our focus on rare events in noisy systems demands that we try to capture mechanism in a probabilistic way, seeking the evolution of the probability distribution for individual (physically interpretable) coordinates. The first example emphasizes that these distributions need not be unimodal; there can be multiple dynamical pathways. Advanced sampling methods like TPS can harvest the reactive ensemble, but that reactive ensemble naturally lets one track the evolution as a function of the time since leaving \mathcal{A} . The reactive ensemble therefore superimposes transitions occurring at stochastically variable times. TPT deconvolves this superposition, allowing us to resolve how the probability distribution over mi-

crostates evolves as a function of reaction progress, q^+ . Our proposed marginalization of the reactive ensemble benefits from being straightforward and simple. Simple, that is, provided the committor can be solved.

Here, we took the direct route to solve for that committor, explicitly representing the vector \mathbf{q}^+ and using sparse linear algebra methods to optimize it. The sparse linear algebra strategy becomes altogether untenable when the state space grows so \mathbf{q}^+ cannot be practically stored in memory. In that case, dimensionality reduction strategies can nevertheless allow the committor to be robustly estimated. Though the number of microstates is astronomically large, the estimated committor might be parameterized by billions or trillions of parameters. For example, tensor train and tensor network approaches can extend sparse linear algebra methods to practically calculate properties of CMEs [79–85], including rare events for large ($\sim 10^{15}$ microstates) reaction-diffusion models [78]. Other approaches using basis expansions [51–54], neural networks [32, 55–61], or tensor networks [62] can even fit high-dimensional committors for the case that \mathbf{x} is continuous. In all of these dimensional reduction strategies, the approximate committor’s parameters (weights in a neural network or elements of a tensor network) are optimized by iterative algorithms. Often, it will be advantageous for those iterative algorithms to learn a committor based from (enhanced) sampling. For example, using trajectories sampled by TPS, one can learn parameters for a deep neural net to estimate the committor [39].

As various strategies for approximating high-dimensional committors develop, it becomes especially important that one can use that committor function to extract information beyond a reaction rate. We expect marginalizations of the reactive ensemble like Eq. (1) to play an important role, but there is a remaining challenge. When the number of microstates becomes astronomical, it is not obvious how one should practically perform the marginalization. Since our vector \mathbf{q}^+ was small enough that we could enumerate it, our $\delta(q(\mathbf{x}) - q)$ of Eq. (1) was implemented by binning each microstate based on its value of the committor and our integral was performed by looping over all microstates. Just because a neural network can estimate the value of q for any given \mathbf{x} , it does not mean that one can simply loop over the microstates to perform the high-dimensional integral. We expect one would approximate Eq. (1) via a Monte Carlo estimate of the high-dimensional integral. Performing those integrals may be more straightforward when committors are approximated by tensor networks since the integration will correspond to traces over physical indices, an operation that is very natural for tensor networks. Even still, the δ function would not be trivial to implement. A candidate is to represent the δ function in a Fourier basis. Though some technical challenges will need to be considered, we expect it will be ultimately be possible to construct the marginals $\rho^{AB}(x_i, q)$ for problems where the curse of dimensionality precludes explicit calculations of \mathbf{q}^+ .

Those marginalizations may also be performed over transformed coordinates. For our examples, we chose to work in a natural set of coordinates captured by (x, y) and $(n_A, n_{A_2}, n_{OA_2}, n_O, n_{OB_2}, n_{B_2}, n_B)$. Our assumption is that these coordinates are simple to physically interpret as positions or counts of particular species. We could alternatively consider correlations between the committor and some transformed physical coordinates. For the first example, one could, for example, perform a rotation to study how distributions for the coordinates $(x + y, x - y)$ evolve with q . In principle, one would not be limited to rotations or linear transformations, but the more complex the transformation, the closer one gets to building q itself. With a sufficiently complex transformation, one loses the physical interpretability and lands on the tautology “the reaction proceeds because the reaction proceeds”. To preserve mechanistic insights like “the reaction proceeds because the distance between x and y shrinks”, it will be important to restrict ourselves to easily interpretable coordinates (the number of species A , the length of a bond, the number of solvent molecules within a radius of a protein, etc.). Then Eq. (1) mixes the benefits of the two types of coordinates, revealing the correlations between simple-to-interpret physical coordinates and the simple-to-interpret concept of reaction progress.

V. ACKNOWLEDGMENTS

We appreciate very useful discussions with Geyao Gu, Emanuele Penocchio, Aaron Dinner, Grant Rotskoff, Spencer Guo, and John Strahan. The material presented in this manuscript is based on work supported by the National Science Foundation under Grant No. 2239867.

Appendix A: Finite truncation

Our ability to generate the marginal reactive ensemble distributions required that we could directly compute \mathbf{q}^+ , something we did in both examples with sparse linear algebra methods. Even with those sparse methods, it is important that one can cap the state space to prevent \mathbf{q}^+ from growing too large. For the GTS model, our imposition of a maximum occupancy, M on non-DNA species served this goal. To test that our cap set at $M = 30$ does not influence the reactive trajectories transitioning between \mathcal{A} and \mathcal{B} , we compared rates calculated by TPT with $M = 30$ to rates computed via forward flux sampling (FFS) [6], as well as the stochastic sampling algorithm (SSA) with no maximum occupancy [15]. From TPT [28],

$$k_{AB} = \sum_{j \in \mathcal{A}, i \notin \mathcal{A} \cup \mathcal{B}} f_{ij}^{AB}, \quad (\text{A1})$$

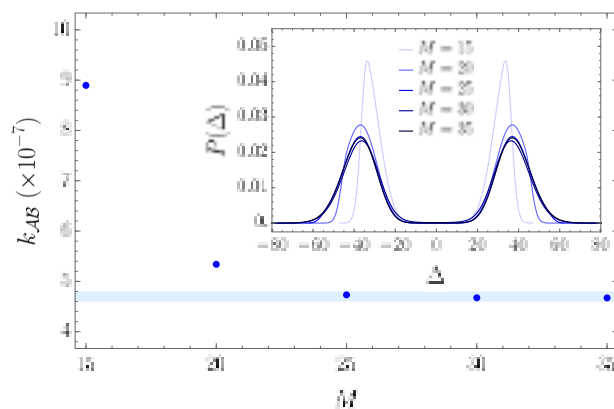


FIG. 4. **Finite truncation convergence.** Switching rates k_{AB} between the two GTS metastable states were computed from Eq. (A1) for various choices of maximum molecule count M . Provided M is sufficiently large, the committor-based rate calculations agree with rates obtained by FFS [6], given by lines with thickness matching the reported standard errors. (Inset) The steady-state distribution for the order parameter Δ also converges for the same sufficiently large M .

where the flux of probability from microstate j to i within the reactive ensemble is

$$f_{ij}^{AB} = \begin{cases} P_i^{AB} W_{ij} & i \neq j, \\ 0, & \text{otherwise} \end{cases}. \quad (\text{A2})$$

A truncation at $M = 30$ was sufficient to yield $k_{AB} = 4.67 \times 10^{-7}$, a rate in excellent agreement with forward flux sampling (FFS) calculations performed on the same model [6] and a brute force rate calculation using 100 SSA trajectories, each of length 10^8 units of time. Fig. 4 shows the convergence of the truncated TPT rates to $k_{AB} = (4.68 \pm 0.05) \times 10^{-7}$, the SSA rate without a truncated maximal occupancy.

The inset of Fig. 4 emphasizes that $M = 30$ was sufficient not only to converge the rate but also to converge distributions. Specifically, we use π to plot the steady-state distribution for Δ . This $P(\Delta)$, which reveals the bimodality for all M , shows that (for the parameters studied) the distribution is only weakly influenced when M exceeds 25.

Appendix B: Avoiding an ill-conditioned backward committor equation

Here, we derive Eq. (8), the linear equation that solves for \mathbf{r} instead of the backward committor \mathbf{q}^- . Observe from Eq. (7) that the reactive ensemble requires a Hadamard product of π , \mathbf{q}^+ , and \mathbf{q}^- . Eq. (6) is an ill-conditioned equation that would solve for \mathbf{q}^- , but we can convert it into a significantly better conditioned equation for \mathbf{r} , the Hadamard product of π and \mathbf{q}^- . To see this conversion, we restrict ourselves to $i, j \in (\mathcal{A} \cup \mathcal{B})^c$ then

substitute $\tilde{U}_{ij} = \tilde{W}_{ji}$ and $\tilde{v}_i = -\sum_{k \in \mathcal{A}} \tilde{W}_{ki}$ into Eq. (6):

$$\sum_j \tilde{W}_{ji} q_j^- = -\sum_{k \in \mathcal{A}} \tilde{W}_{ki}. \quad (\text{B1})$$

Rewriting this equation in terms of the time-forward matrix W , we have

$$\frac{1}{\pi_i} \sum_j W_{ij} \pi_j q_j^- = -\frac{1}{\pi_i} \sum_{k \in \mathcal{A}} W_{ik} \pi_k. \quad (\text{B2})$$

For Eq. (B2) hold for all i , we therefore require

$$\sum_j W_{ij} \pi_j = -\sum_{k \in \mathcal{A}} W_{ik} \pi_k. \quad (\text{B3})$$

Finally, recalling that U and W are transposes of each other within the $(\mathcal{A} \cup \mathcal{B})^c$ region, the expression simplifies to Eq. 8:

$$U^T \mathbf{r} = \mathbf{s}. \quad (\text{B4})$$

Appendix C: Doi-Peliti construction of W

A microstate of the GTS is given by $(n_A, n_{A_2}, n_{OA_2}, n_O, n_{OB_2}, n_{B_2}, n_B)$. Recognizing that the DNA exists in the OA_2 , O , or OB_2 state, we equivalently define $\mathbf{n} = (n_A, n_{A_2}, n_\lambda, n_{B_2}, n_B)$, where $n_\lambda = 0, 1$, and 2 correspond to the OA_2 , O , and OB_2 states, respectively. The vector of probabilities of each microstate, \mathbf{p}_t , evolves according to the master equation

$$\frac{d\mathbf{p}_t}{dt} = W\mathbf{p}_t, \quad (\text{C1})$$

where W is a rate operator constructed from the 14 reactions of Fig. 2. Writing that W in matrix form can be an accounting headache that requires one to enumerate the microstates. Instead, it can be convenient to write both \mathbf{p}_t and W in a tensor-product form that isolates the action of each reaction on the occupation numbers of the chemical species. Here, we sketch the framework for constructing W in terms of operators that raise and lower $n_A, n_{A_2}, n_\lambda, n_{B_2}$, and n_B . Readers interested in more algebraic details are referred to the appendices of Ref. [78].

The aim is to write each reaction's contribution to the rate operator W in a tensor-product form:

$$O_A \otimes O_{A_2} \otimes O_\lambda \otimes O_{B_2} \otimes O_B, \quad (\text{C2})$$

where each O_γ is an operator that acts on a local state space spanned by the possible states of $|n_\gamma\rangle$. Since A, A_2, B , and B_2 have an occupancy number between 0 and M , their local state spaces are spanned by orthonormal basis vectors $|0\rangle, |1\rangle \dots |M\rangle$, which means that the operators acting on their local state spaces are merely $(M+1) \times (M+1)$ matrices. Operators on the λ space are even smaller—they are simply 3×3 matrices. The states of the many-body system are spanned by the tensor product

states $|\mathbf{n}\rangle = |n_A\rangle \otimes |n_{A_2}\rangle \otimes |n_\lambda\rangle \otimes |n_{B_2}\rangle \otimes |n_B\rangle$, which are also orthonormal. We write a probability distribution over microstates as a superposition of the many-body states:

$$|p_t\rangle = \sum_{\mathbf{n}} p_t(\mathbf{n}) |\mathbf{n}\rangle. \quad (\text{C3})$$

Eq. (C3) is the tensor-product form of what we previously called \mathbf{p}_t . Inspecting how each reaction impacts $p_t(\mathbf{n})$, we are now in a position to build the tensor-product form of W .

To gain an intuition about how a chemical reaction converts into the set of local operators, it is useful to explicitly consider the first reaction of Fig. 2, $2A \xrightarrow{k_1} A_2$. The action of this reaction is to decrease n_A by two and to increase n_{A_2} by one, so it is useful to define a raising operator x_γ^\dagger and a corresponding lowering operator x_γ that act on species gamma. Taking into account that species γ has a maximum occupancy of M_γ , these operators are defined by

$$\begin{aligned} x_\gamma |n_\gamma\rangle &= \begin{cases} n_\gamma |n_\gamma - 1\rangle, & 0 < n_\gamma \leq M_\gamma, \\ 0, & \text{otherwise,} \end{cases} \\ x_\gamma^\dagger |n_\gamma\rangle &= \begin{cases} |n_\gamma + 1\rangle, & 0 \leq n_\gamma < M_\gamma - 1, \\ 0, & \text{otherwise.} \end{cases} \end{aligned} \quad (\text{C4})$$

In matrix form,

$$x = \begin{pmatrix} 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 2 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & M \\ 0 & 0 & 0 & \dots & 0 & 0 \end{pmatrix} \text{ and } x^\dagger = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 & 0 \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 1 & 0 \end{pmatrix}. \quad (\text{C5})$$

Therefore, one might guess that reaction 1 contributes to W a term of the form $x_A^2 \otimes x_{A_2}^\dagger \otimes \mathbb{I}_\lambda \otimes \mathbb{I}_{B_2} \otimes \mathbb{I}_B$, a guess that involves lowering A twice, raising A_2 once, and acting on the other species with the identity \mathbb{I} to leave them unchanged. That guess correctly anticipates the off-diagonal components of W , but to conserve probability, there is an additional negative element along the diagonal of W . That negative term is especially clear in the gain-loss CME for the first reaction:

$$\begin{aligned} \frac{dp_t(\mathbf{n})}{dt} &= k_1 \left[(n_A + 2)(n_A + 1)p_t(n_A + 2, n_{A_2} - 1) \right. \\ &\quad \left. - n_A(n_A - 1)p_t(n_A, n_{A_2}) \right] \end{aligned} \quad (\text{C6})$$

By summing both sides of Eq. (C6) over microstates (i.e., $\sum_{\mathbf{n}} \dots |\mathbf{n}\rangle$) and by judiciously replacing terms like n_A by their number operator representation $a^\dagger a$, the ac-

tion of reaction 1 can be expressed as

$$\frac{d|p_t\rangle}{dt} = k_1 \left(x_A^2 \otimes x_{A_2}^\dagger \otimes \mathbb{I}_\lambda \otimes \mathbb{I}_{B_2} \otimes \mathbb{I}_B \right. \\ \left. - x_A^{\dagger 2} x_{A_2}^2 \otimes y_{A_2} \otimes \mathbb{I}_\lambda \otimes \mathbb{I}_{B_2} \otimes \mathbb{I}_B \right) |p_t\rangle, \quad (C7)$$

where

$$y = \mathbb{I} - |M\rangle \langle M| - |M-1\rangle \langle M-1| = \begin{pmatrix} 1 & \cdots & 0 & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & \cdots & 1 & 0 & 0 \\ 0 & \cdots & 0 & 0 & 0 \\ 0 & \cdots & 0 & 0 & 0 \end{pmatrix} \quad (C8)$$

adjusts the probability conserving diagonal element to accommodate for the fact that $n_{A_2} = M+2 \rightarrow M$ transitions have been removed by the truncation [78]. To compress the notation, it is customary to suppress the identity operators and the tensor product symbols, writing reaction 1's contribution to W as simply

$$W^{2A \rightarrow A_2} = k_1 \left(x_A^2 x_{A_2}^\dagger - x_A^{\dagger 2} x_{A_2}^2 y_{A_2} \right). \quad (C9)$$

Similar procedures can be carried out for the other 13 reactions.

$$\begin{aligned} W^{2A \rightarrow A_2} &= k_1 \left(x_A^2 x_{A_2}^\dagger - x_A^{\dagger 2} x_{A_2}^2 y_{A_2} \right) \\ W^{A_2 \rightarrow 2A} &= k_2 \left(x_A^{\dagger 2} x_{A_2} - z_A x_{A_2}^\dagger x_{A_2} \right) \\ W^{O+A_2 \rightarrow OA_2} &= k_3 \left(x_{A_2} a_\lambda^\dagger - x_{A_2}^\dagger x_{A_2} \omega_\lambda \right) \\ W^{OA_2 \rightarrow O+A_2} &= k_4 \left(x_{A_2}^\dagger a_\lambda - z_{A_2} \alpha_\lambda \right) \\ W^{O \rightarrow O+A} &= k_5 \left(x_A^\dagger \omega_\lambda - z_A \omega_\lambda \right) \\ W^{OA_2 \rightarrow OA_2+A} &= k_6 \left(x_A^\dagger \alpha_\lambda - z_A \alpha_\lambda \right) \\ W^{A \rightarrow \emptyset} &= k_7 \left(x_A - x_A^\dagger x_A \right) \end{aligned}$$

Here, we have introduced this tensor-product form as a convenient way to construct W for sparse matrix operations, but we note that it is also the starting point

For compactness, it is useful to additionally define

$$z = \mathbb{I} - |M\rangle \langle M| = \begin{pmatrix} 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 0 \end{pmatrix} \quad (C10)$$

as well as a set of 3×3 operators acting on the λ space:

$$\begin{aligned} a^\dagger &= \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad b^\dagger = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}, \\ a &= \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad b = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}, \\ \alpha &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad \beta = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \\ \omega &= \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}. \end{aligned} \quad (C11)$$

The operators a^\dagger and a respectively create and destroy OA_2 from O while b^\dagger and b respectively create and destroy OB_2 . The final three operators, α , β , and ω detect the occupancy of the OA_2 , OB_2 , and O states, respectively. Having defined all the necessary local operators, we finally write down the contribution to W from each of the 14 reactions:

$$\begin{aligned} W^{2B \rightarrow B_2} &= k_1 \left(x_{B_2}^\dagger x_B^2 - y_{B_2} x_B^{\dagger 2} x_B^2 \right) \\ W^{B_2 \rightarrow 2B} &= k_2 \left(x_{B_2} x_B^{\dagger 2} - x_{B_2}^\dagger x_{B_2} z_B \right) \\ W^{O+B_2 \rightarrow OB_2} &= k_3 \left(b_\lambda^\dagger x_{B_2} - \omega_\lambda x_{B_2}^\dagger x_{B_2} \right) \\ W^{OB_2 \rightarrow O+B_2} &= k_4 \left(b_\lambda x_{B_2}^\dagger - \beta_\lambda z_{B_2} \right) \\ W^{O \rightarrow O+B} &= k_5 \left(\omega_\lambda x_B^\dagger - \omega_\lambda z_B \right) \\ W^{OB_2 \rightarrow OB_2+B} &= k_6 \left(\beta_\lambda x_B^\dagger - \beta_\lambda z_B \right) \\ W^{B \rightarrow \emptyset} &= k_7 \left(x_B - x_B^\dagger x_B \right). \end{aligned}$$

for employing tensor network methods. Those methods promise to make these committor calculations practical for even larger systems.

[1] K. Svoboda, C. F. Schmidt, B. J. Schnapp, and S. M. Block, Direct observation of kinesin stepping by optical

trapping interferometry, *Nature* **365**, 721 (1993).

- [2] H. H. McAdams and A. Arkin, Stochastic mechanisms in gene expression, *Proc. Natl. Acad. Sci. USA* **94**, 814 (1997).
- [3] X. Zhuang, L. E. Bartley, H. P. Babcock, R. Russell, T. Ha, D. Herschlag, and S. Chu, A single-molecule study of RNA catalysis and folding, *Science* **288**, 2048 (2000).
- [4] W. A. Eaton, V. Muñoz, S. J. Hagen, G. S. Jas, L. J. Lapidus, E. R. Henry, and J. Hofrichter, Fast kinetics and mechanisms in protein folding, *Annu. Rev. Biophys. Biomol. Struct.* **29**, 327 (2000).
- [5] I. Golding, J. Paulsson, S. M. Zawilski, and E. C. Cox, Real-time kinetics of gene activity in individual bacteria, *Cell* **123**, 1025 (2005).
- [6] R. J. Allen, P. B. Warren, and P. R. ten Wolde, Sampling rare switching events in biochemical networks, *Phys. Rev. Lett.* **94**, 018104 (2005).
- [7] A. Eldar and M. B. Elowitz, Functional roles for noise in genetic circuits, *Nature* **467**, 167 (2010).
- [8] D. T. Gillespie, A rigorous derivation of the chemical master equation, *Physica A* **188**, 404 (1992).
- [9] N. van Kampen, *Stochastic Processes in Physics and Chemistry*, Vol. 1 (Elsevier, 1992).
- [10] J. D. Chodera and F. Noé, Markov state models of biomolecular conformational dynamics, *Curr. Opin. Struct. Biol.* **25**, 135 (2014).
- [11] B. E. Husic and V. S. Pande, Markov state models: From an art to a science, *J. Am. Chem. Soc.* **140**, 2386 (2018).
- [12] P. Thomas, N. Popović, and R. Grima, Phenotypic switching in gene regulatory networks, *Proc. Natl. Acad. Sci. USA* **111**, 6994 (2014).
- [13] G. Graciun, Y. Tang, and M. Feinberg, Understanding bistability in complex enzyme-driven reaction networks, *Proc. Natl. Acad. Sci. U.S.A.* **103**, 8697 (2006).
- [14] E. Ozbudak, M. Thattai, H. N. Lim, B. I. Shraiman, and A. van Oudenaarden, Multistability in the lactose utilization network of *Escherichia coli*, *Nature* **427**, 737 (2004).
- [15] D. T. Gillespie, A. Hellander, and L. R. Petzold, Perspective: Stochastic algorithms for chemical kinetics, *J. Chem. Phys.* **138**, 170901 (2013).
- [16] C. Yang, D. T. Gillespie, and L. R. Petzold, Efficient step size selection for the tau-leaping simulation method, *J. Chem. Phys.* **124**, 044104 (2006).
- [17] P. Hänggi, P. Talkner, and M. Borkovec, Reaction-rate theory: Fifty years after Kramers, *Rev. Mod. Phys.* **62**, 251 (1990).
- [18] C. Dellago, P. G. Bolhuis, F. S. Csajka, and D. Chandler, Transition path sampling and the calculation of rate constants, *J. Chem. Phys.* **108**, 1964 (1998).
- [19] P. G. Bolhuis, D. Chandler, C. Dellago, and P. L. Geissler, Transition path sampling: Throwing ropes over rough mountain passes, in the dark, *Annu. Rev. Phys. Chem.* **53**, 291 (2002).
- [20] G. E. Crooks and D. Chandler, Efficient transition path sampling for nonequilibrium stochastic dynamics, *Phys. Rev. E* **64**, 026109 (2001).
- [21] T. R. Gingrich and P. L. Geissler, Preserving correlations between trajectories for efficient path sampling, *J. Chem. Phys.* **142**, 234104 (2015).
- [22] R. J. Allen, D. Frenkel, and P. R. ten Wolde, Simulating rare events in equilibrium or nonequilibrium stochastic systems, *J. Chem. Phys.* **124**, 024102 (2006).
- [23] R. J. Allen, C. Valeriani, and P. R. ten Wolde, Forward flux sampling for rare event simulations, *J. Condens. Matter Phys.* **21**, 463102 (2009).
- [24] W. E and E. Vanden-Eijnden, Transition-path theory and path-finding algorithms for the study of rare events, *Annu. Rev. Phys. Chem.* **61**, 391 (2010).
- [25] J. Lu and E. Vanden-Eijnden, Exact dynamical coarse-graining without time-scale separation, *J. Chem. Phys.* **141**, 044109 (2014).
- [26] W. E and E. Vanden-Eijnden, Towards a theory of transition paths, *J. Stat. Phys.* **123**, 503 (2006).
- [27] P. Metzner, C. Schütte, and E. Vanden-Eijnden, Illustration of transition path theory on a collection of simple examples, *J. Chem. Phys.* **125**, 084110 (2006).
- [28] P. Metzner, Schütte, and E. Vanden-Eijnden, Transition path theory for Markov jump processes, *Multiscale Model. Simul.* **7**, 1192 (2009).
- [29] R. Du, V. S. Pande, A. Y. Grosberg, T. Tanaka, and E. S. Shakhnovich, On the transition coordinate for protein folding, *J. Chem. Phys.* **108**, 334 (1998).
- [30] P. G. Bolhuis, C. Dellago, and D. Chandler, Reaction coordinates of biomolecular isomerization, *Proc. Natl. Acad. Sci. USA* **97**, 5877 (2000).
- [31] Y. M. Rhee and V. S. Pande, One-dimensional reaction coordinate and the corresponding potential of mean force from commitment probability distribution, *J. Phys. Chem. B* **109**, 6780 (2005).
- [32] A. Ma and A. R. Dinner, Automatic method for identifying reaction coordinates in complex systems, *J. Phys. Chem. B* **109**, 6769 (2005).
- [33] W. Lechner, J. Rogal, J. Juraszek, B. Ensing, and P. G. Bolhuis, Nonlinear reaction coordinate analysis in the reweighted path ensemble, *J. Chem. Phys.* **133**, 174110 (2010).
- [34] M. A. Rohrdanz, W. Zheng, and C. Clementi, Discovering mountain passes via torchlight: Methods for the definition of reaction coordinates and pathways in complex macromolecular reactions, *Annu. Rev. Phys. Chem.* **64**, 295 (2013).
- [35] K. Neupane, A. P. Manuel, and M. T. Woodside, Protein folding trajectories can be described quantitatively by one-dimensional diffusion over measured energy landscapes, *Nat. Phys.* **12**, 700 (2016).
- [36] R. Elber, J. M. Bello-Rivas, P. Ma, A. E. Cardenas, and A. Fathizadeh, Calculating iso-committor surfaces as optimal reaction coordinates with milestoning, *Entropy* **19**, 219 (2017).
- [37] Z. He, C. Chipot, and B. Roux, Committor-consistent variational string method, *J. Phys. Chem. Lett.* **13**, 9263 (2022).
- [38] B. Roux, Transition rate theory, spectral analysis, and reactive paths, *J. Chem. Phys.* **156**, 134111 (2022).
- [39] H. Jung, R. Covino, A. Arjun, C. Leitold, C. Dellago, P. G. Bolhuis, and G. Hummer, Machine-guided path sampling to discover mechanisms of molecular self-organization, *Nat. Comput. Sci.* **3**, 334 (2023).
- [40] L. Evans, M. K. Cameron, and P. Tiwary, Computing committors in collective variables via Mahalanobis diffusion maps, *Appl. Comput. Harmon. Anal.* **64**, 62 (2023).
- [41] P. L. Geissler, C. Dellago, and D. Chandler, Kinetic pathways of ion pair dissociation in water, *J. Phys. Chem. B* **103**, 3706 (1999).
- [42] G. Hummer, From transition paths to transition states and rate coefficients, *J. Chem. Phys.* **120**, 516 (2004).
- [43] R. B. Best and G. Hummer, Reaction coordinates and rates from transition paths, *Proc. Natl. Acad. Sci. USA*

- [43] 102, 6732 (2005).
- [44] P. V. Banushkina and S. V. Krivov, Optimal reaction coordinates, *WIREs Comput. Mol. Sci.* **6**, 748 (2016).
- [45] B. Peters, Reaction coordinates and mechanistic hypothesis tests, *Annu. Rev. Phys. Chem.* **67**, 669 (2016).
- [46] W. Zhang, C. Hartmann, and C. Schütte, Effective dynamics along given reaction coordinates, and reaction rate theory, *Faraday Discuss.* **195**, 365 (2016).
- [47] S. V. Krivov, Protein folding free energy landscape along the committor-the optimal folding coordinate, *J. Chem. Theory Comput.* **14**, 3418 (2018).
- [48] S. Wu, H. Li, and A. Ma, A rigorous method for identifying a one-dimensional reaction coordinate in complex molecules, *J. Chem. Theory Comput.* **18**, 2836 (2022).
- [49] L. Mouaffac, K. Palacio-Rodriguez, and F. Pietrucci, Optimal reaction coordinates and kinetic rates from the projected dynamics of transition paths, *J. Chem. Theory Comput.* **19**, 5701 (2023).
- [50] H. Chen, B. Roux, and C. Chipot, Discovering reaction pathways, slow variables, and committor probabilities with machine learning, *J. Chem. Theory Comput.* (2023).
- [51] E. H. Thiede, D. Giannakis, A. R. Dinner, and J. Weare, Galerkin approximation of dynamical quantities using trajectory data, *J. Chem. Phys.* **150**, 244111 (2019).
- [52] J. Strahan, A. Antoszewski, C. Lorpaiboon, B. P. Vani, J. Weare, and A. R. Dinner, Long-time-scale predictions from short-trajectory data: A benchmark analysis of the Trp-cage miniprotein, *J. Chem. Theory Comput.* **17**, 2948–2963 (2021).
- [53] D. Aristoff, M. Johnson, G. Simpson, and R. Webber, The fast committor machine: Interpretable prediction with kernels, *arXiv:2405.10410* 10.48550/arXiv.2405.10410 (2024).
- [54] S. C. Guo, R. Shen, B. Roux, and A. R. Dinner, Dynamics of activation in the voltage-sensing domain of *Ciona intestinalis* phosphatase Ci-VSP, *Nat. Commun.* **15**, 1408 (2024).
- [55] G. M. Rotskoff, A. R. Mitchell, and E. Vanden-Eijnden, Active importance sampling for variational objectives dominated by rare events: Consequences for optimization and generalization, in *Mathematical and Scientific Machine Learning* (PMLR, 2022) pp. 757–780.
- [56] J. Strahan, J. Finkel, A. R. Dinner, and J. Weare, Predicting rare events using neural networks and short-trajectory data, *J. Comput. Phys.* **488**, 112152 (2023).
- [57] Y. Khoo, J. Lu, and L. Ying, Solving for high-dimensional committor functions using artificial neural networks, *Res. Math. Sci.* **6**, 1 (2018).
- [58] M. R. Hasyim, C. H. Batton, and K. K. Mandadapu, Supervised learning and the finite-temperature string method for computing committor functions and reaction rates, *J. Chem. Phys.* **157**, 184111 (2022).
- [59] Q. Li, B. Lin, and W. Ren, Computing committor functions for the study of rare events using deep learning, *J. Chem. Phys.* **151**, 054112 (2019).
- [60] V. Jacques-Dumas, R. M. van Westen, F. Bouchet, and H. A. Dijkstra, Data-driven methods to estimate the committor function in conceptual ocean models, *Nonlinear Process. Geophys.* **30**, 195 (2023).
- [61] B. Lin and W. Ren, Deep learning method for computing committor functions with adaptive sampling, *arXiv:2404.06206* 10.48550/arXiv.2404.06206 (2024).
- [62] Y. Chen, J. Hoskins, Y. Khoo, and M. Lindsey, Commit-tor functions via tensor networks, *J. Comput. Phys.* **472**, 111646 (2023).
- [63] F. Noé, C. Schütte, E. Vanden-Eijnden, L. Reich, and T. R. Weiki, Constructing the equilibrium ensemble of folding pathways from short off-equilibrium simulations, *Proc. Natl Acad. Sci. U.S.A.* **106**, 19011 (2009).
- [64] T. J. Lane, G. R. Bowman, K. Beauchamp, V. A. Voelz, and V. S. Pande, Markov state model reveals folding and functional dynamics in ultra-long MD trajectories, *J. Am. Chem. Soc.* **133**, 218413 (2011).
- [65] J. Finkel, R. J. Webber, E. P. Gerber, D. S. Abbot, and J. Weare, Data-driven transition path analysis yields a statistical understanding of sudden stratospheric warming events in an idealized model, *JAS* **80**, 519 (2023).
- [66] M. Cameron and E. Vanden-Eijnden, Flows in complex networks: Theory, algorithms, and application to Lennard-Jones cluster rearrangement, *J. Stat. Phys.* **156**, 427 (2014).
- [67] M. D. Louwerse and D. A. Sivak, Information thermodynamics of the transition-path ensemble, *Phys. Rev. Lett.* **128**, 170602 (2022).
- [68] V. A. Voelz, G. R. Bowman, and K. Beauchamp, Molecular simulation of ab initio protein folding for a millisecond folder NTL9(1–39), *J. Am. Chem. Soc.* **132**, 1526 (2010).
- [69] D. Shuklar, Y. Meng, B. Roux, and V. S. Pande, Activation pathway of src kinase reveals intermediate states as targets for drug design, *J. Am. Chem. Soc.* **132**, 1526 (2010).
- [70] E. Guamera and E. Vanden-Eijnden, Optimized Markov state models for metastable systems, *J. Chem. Phys.* **145**, 024102 (2016).
- [71] Y. Liu, D. P. Hickey, S. D. Minter, A. Dickson, and S. C. Barton, Markov-state transition path analysis of electrostatic channeling, *J. Phys. Chem. C* **123**, 15284 (2019).
- [72] R. Banerjee and R. I. Cukier, Transition paths of Met-Enkephalin from Markov state modeling of a molecular dynamics trajectory, *J. Phys. Chem. B* **118**, 2883 (2014).
- [73] B. P. Vani, J. Weare, and A. R. Dinner, Computing transition path theory quantities with trajectory stratification, *J. Chem. Phys.* **157**, 034106 (2022).
- [74] B. W. Silverman, *Density estimation*, Vol. 1 (London: Chapman and Hall, 1986).
- [75] N. E. Strand, R.-S. Fu, and T. R. Gingrich, Current inversion in a periodically driven two-dimensional Brownian ratchet, *Phys. Rev. E* **102**, 012141 (2020).
- [76] M. Doi, Stochastic theory of diffusion-controlled reaction, *J. Phys. A* **9**, 1479 (1976).
- [77] L. Peliti, Path integral approach to birth-death processes on a lattice, *J. Phys.* **46**, 1469 (1985).
- [78] S. B. Nicholson and T. R. Gingrich, Quantifying rare events in stochastic reaction-diffusion dynamics using tensor networks, *Phys. Rev. X* **13**, 041006 (2023).
- [79] M. Hegland and J. Garcke, On the numerical solution of the chemical master equation with sums of rank one tensors, *ANZIAM J.* **52**, C628 (2010).
- [80] V. Kazeev, M. Khammash, M. Nip, and C. Schwab, Direct solution of the chemical master equation using quantized tensor trains, *PLoS Computat. Bio.* **10**, e1003359 (2014).
- [81] S. Liao, T. Vejchodský, and R. Erban, Tensor methods for parameter estimation and bifurcation analysis of stochastic reaction networks, *J. R. Soc. Interface* **12**, 20150233 (2015).

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0232705

- [82] S. Dolgov and B. Khoromskij, Simultaneous state-time approximation of the chemical master equation using tensor product formats, *Numer. Linear Algebra Appl.* **22**, 197 (2015).
- [83] P. Gelß, S. Matera, and C. Schütte, Solving the master equation without kinetic Monte Carlo: Tensor train approximations for a CO oxidation model, *J. Comput. Phys.* **314**, 489 (2016).
- [84] H. D. Vo and R. B. Sidje, An adaptive solution to the chemical master equation using tensors, *J. Chem. Phys.* **147**, 044102 (2017).
- [85] I. G. Ion, C. Wildner, D. Loukrezis, H. Koeppel, and H. De Gerssem, Tensor-train approximation of the chemical master equation and its application for parameter inference, *J. Chem. Phys.* **155**, 034102 (2021).