

Gage Bonner ■ (D); F. J. Beron-Vera (D); M. J. Olascoaga (D)



Chaos 33, 063141 (2023)

https://doi.org/10.1063/5.0144706





CrossMark





Improving the stability of temporal statistics in transition path theory with sparse data 0

Cite as: Chaos 33, 063141 (2023); doi: 10.1063/5.0144706 Submitted: 31 January 2023 · Accepted: 26 May 2023 ·











Published Online: 15 June 2023



Gage Bonner, 1, a) F. J. Beron-Vera, 1, b) and M. J. Olascoaga^{2,c)}



AFFILIATIONS

Department of Atmospheric Sciences, Rosenstiel School of Marine, Atmospheric, and Earth Science, University of Miami, Miami, Florida 33149, USA

²Department of Ocean Sciences, Rosenstiel School of Marine, Atmospheric, and Earth Science, University of Miami, Miami, Florida 33149, USA

^{a)}Author to whom correspondence should be addressed: gbonner@miami.edu

b) Electronic mail: fberon@miami.edu

c)Electronic mail: jolascoaga@miami.edu

ABSTRACT

Ulam's method is a popular discretization scheme for stochastic operators that involves the construction of a transition probability matrix controlling a Markov chain on a set of cells covering some domain. We consider an application to satellite-tracked undrogued surface-ocean drifting buoy trajectories obtained from the National Oceanic and Atmospheric Administration Global Drifter Program dataset. Motivated by the motion of Sargassum in the tropical Atlantic, we apply Transition Path Theory (TPT) to drifters originating off the west coast of Africa to the Gulf of Mexico. We find that the most common case of a regular covering by equal longitude-latitude side cells can lead to a large instability in the computed transition times as a function of the number of cells used. We propose a different covering based on a clustering of the trajectory data that is stable against the number of cells in the covering. We also propose a generalization of the standard transition time statistic of TPT that can be used to construct a partition of the domain of interest into weakly dynamically connected regions.

Published under an exclusive license by AIP Publishing. https://doi.org/10.1063/5.0144706

Transition Path Theory (TPT) provides a rigorous statistical characterization of the ensemble of trajectories connecting directly, i.e., without detours, two disconnected (sets of) states in a Markov chain, a stochastic process that undergoes transitions from one state to another with probability depending on the state attained in the previous step. Markov chains can be constructed using trajectory data via counting of transitions between cells covering the domain spanned by trajectories. With sparse trajectory data, the use of regular cells is observed to result in unstable estimates of the total duration of transition paths. Using Voronoi cells resulting from k-means clustering of the trajectory data, we obtain stable estimates of this TPT statistic, which is generalized to frame the remaining duration of transition paths. This remaining duration is a new TPT statistic suitable for investigating connectivity.

I. INTRODUCTION

Sargassum is a pelagic seaweed that plays a crucial role in the ecosystem of the Sargasso Sea and surrounding areas of the North Atlantic.1 Large rafts of the seaweed drift through the Caribbean and into the Gulf of Mexico before being circulated into the Sargasso Sea by the Gulf Stream where it is replenished yearly.² These Sargassum clumps provide a habitat to a diverse contingent of invertebrate, fish, and other fauna far away from land. In addition, the Sargasso Sea contributes to approximately 7% of the global net biological carbon pump due to the abundance of Sargassum and the community of organisms it houses.3 In 2011, islands in the Caribbean Sea and beaches in South Florida were inundated with abnormally large quantities of Sargassum. Since then, waves of Sargassum have been reported regularly in these locations as well as in western Africa and northern Brazil. Although beached Sargassum can add nutrients to

coastal soils,³ it also creates offensive smells and can result in the destruction of some habitats.⁵ Large scale cleanup of beaches costs millions annually and can negatively impact tourism in the affected regions.

The study of the transport of Sargassum across the ocean has attracted much interest. Satellite-tracked drifter trajectory data from the National Oceanic and Atmospheric Administration (NOAA) Global Drifter Program (GDP)⁶ have been used to infer the evolution of the density of Sargassum. Since Sargassum tends to remain on the surface of the ocean, windage effects must be taken into account. Drifters are placed into the ocean with a drogue, a kind of anchor that helps the drifter move with the surface currents, by resisting wind slippage and wave-induced drift. Drogues often detach themselves after some time, leaving a drifter that is more susceptible to wind and wave effects. In Beron-Vera et al.,7 motivated by this consideration, it is demonstrated that the motion of undrogued drifters tracks the actual satellite-inferred density of Sargassum more closely than the drogued counterparts. To do this, the North Atlantic is discretized into a large number of small cells, which define the states of a Markov chain whose transition probability matrix is constructed based on the initial and final locations of drifter trajectory data on a certain time interval. Provided certain technical conditions are satisfied by this Markov chain, Transition Path Theory (TPT)^{8,9} can be applied to identify bottlenecks and fluxes between a source and a target state. In Beron-Vera et al., by taking the source to be a single cell off the coast of West Africa, and the target to be the Gulf of Mexico, two prominent paths taken by drifters are revealed. The first is a "direct" path along the Great Atlantic Sargassum Belt;10 the shape of this path is in agreement with the satellite-derived density of Sargassum in this region. The second is an "indirect" southern path whereby drifters circulate toward the Gulf of Guinea before eventually traveling westward along the coast of northern Brazil into the Caribbean.1

The time taken to transition between the source and the target is another statistic that can be computed using TPT. It was noticed that the method of Beron-Vera et al.7 provides a transition time, which is highly sensitive to the number of boxes chosen to cover the domain, creating some distrust in the results. This motivated the development of a new kind of covering, which leads to transition times that are stable as a function of the number of cells in the partition. Briefly, this involves clustering the data and generating a covering based on the boundaries of the clusters. Thus, requesting a finer grid tends to result in the division of larger cells while leaving distant cells unchanged. Since we take the stability of the transition time as a metric for the trustworthiness of the application of TPT, we review this statistic by proposing a more general one. This new statistic has the properties that it (1) gives the standard transition time as a special case (so long as the source is a single cell) and (2) provides means for partitioning the flow domain to investigate connectivity.

The remainder of this paper is organized as follows. In Sec. II, we review the theoretical framework of the discretization scheme we apply to the trajectory data in our domain. In Sec. III, we apply this discretization scheme to obtain a Markov chain suitable for the application of TPT. We compute transition times and other statistics for two standard kinds of coverings based on regular grids of squares and hexagons to understand their shortcomings. We then propose a different discretization scheme based on the k-means clustering

algorithm which is then shown to be significantly more stable. We also investigate the effect of the transition time step through which the trajectory data are temporally "sliced" for both regular coverings and our new covering. In Sec. IV, we introduce our generalized transition time and demonstrate how it can be used to obtain a partition of our domain into weakly dynamically connected regions. The proof that our generalized transition time reduces to the standard transition time of TPT in the appropriate limit is given in the Appendix. Finally, Sec. V summarizes our results and conclusions.

II. BACKGROUND

A. Trajectory discretization

We consider data sets consisting of a series of J disconnected trajectories $x_1(t), x_2(t), \ldots, x_J(t)$ in \mathbb{X} , where \mathbb{X} is a subset of the two-sphere. Each trajectory $x_i(t)$ consists of a number of observations regularly spaced in time by Δt units. We suppose that each trajectory is generated by the same underlying nondeterministic dynamical map \mathcal{L} , which takes elements of \mathbb{X} to \mathbb{X} -valued random variables on the appropriate probability space equipped with Lebesgue measure m. Let $f \in L^1(\mathbb{X})$ be almost-every non-negative and normalized such that $||f||_{L^1} = 1$. If \mathcal{L} has a stochastic kernel $K(x,y): \mathbb{X} \times \mathbb{X} \to \mathbb{R}^+$ such that $\mathcal{L}(x) \sim K(x,\cdot)$, where $\int_{\mathbb{X}} K(\cdot,y) \, dm(y) = 1$, then we can define the Perron–Frobenius operator, also known as a transfer operator, $\mathcal{P}: L^1(\mathbb{X}) \to L^1(\mathbb{X})$ as $L^1(\mathbb{X}) \to L^1(\mathbb{X})$

$$\mathcal{P}f(y) = \int_{\mathbb{T}} f(x)K(x,y) \, \mathrm{d}m(x). \tag{1}$$

The Perron–Frobenius operator describes how an initial distribution is pushed forward by the underlying dynamics. We can study the action of $\mathcal L$ numerically by discretizing the Perron–Frobenius operator and using the known trajectory data.

The most widely used discretization scheme is Ulam's method.^{13,14} Let $\{B_1, \ldots, B_N\}$ be a partition of $\mathbb X$ into disjoint sets and let $\mathbf{1}_{\mathbb B}(x)$ be the indicator function on the set $\mathbb B$, which gives 1 when $x \in \mathbb B$ and 0 otherwise. Ulam's method can be interpreted as a Galerkin projection¹⁵ of $\mathcal P$ onto the subspace spanned by $\{\mathbf{1}_{B_1}, \ldots, \mathbf{1}_{B_N}\}$. By choosing basis functions $\{m(B_i)^{-1}\mathbf{1}_{B_i}(x)\}$, we have that the discretization of $\mathcal P$ is an N-dimensional linear operator $\mathbf P$ given by a matrix (P_{ij}) such that $P_{i,j}$ is the transition probability from B_i to B_j . ^{16,17} Since $\mathcal P$ acts on $\{m(B_i)^{-1}\mathbf{1}_{B_i}(x)\}$, $P_{i,j}$, by Eq. (1), we have

$$P_{i,j} = \int_{B_j} \mathcal{P} \frac{\mathbf{1}_{B_i}(x)}{m(B_i)} \, \mathrm{d}\mu = \frac{1}{m(B_i)} \int_{B_i} \int_{B_j} K(x, y) \, \mathrm{d}m(x) \, \mathrm{d}m(y). \tag{2}$$

The matrix **P** is a row-stochastic transition probability matrix, which is the discretized analog of K(x, y). Note that the factor $m(B_i)^{-1}$ in the choice of basis functions is what ensures that **P** is (row) stochasticized. For computational purposes, we approximate Eq. (2) in terms of the trajectory data as¹⁶

$$P_{ij} \approx \frac{\sum_{\ell=1}^{J} \sum_{t} \mathbf{1}_{B_i}(x_{\ell}(t)) \mathbf{1}_{B_j}(x_{\ell}(t+T))}{\sum_{\ell=1}^{J} \sum_{t} \mathbf{1}_{B_i}(x_{\ell}(t))},$$
 (3)

where T is some multiple of Δt . This approach has been used in numerous applications. ^{18–28} The transition probability matrix **P** defines a Markov chain on N states such that the ith state is thought

of as a delta distribution of mass located at the center of B_i . In a practical setting, T must be chosen large enough such that the Markov property holds to suitable precision. In summary, the procedure for the translation of trajectory data into a transition probability matrix involves two main degrees of freedom: (1) a choice of covering of the computational domain by disjoint boxes and (2) a choice of T. We return to these issues in Secs. III A and III C, respectively.

B. Transition path theory

We summarize the key results of Transition Path Theory (TPT) here; details can be found in a series of works. 8,9,29,30 We use $\Pr(\cdot)$ and $\operatorname{Ex}[\cdot]$ to indicate probabilities and expectations, respectively. To begin, we consider a discrete Markov chain $(X_n)_{n\in\mathbb{Z}}$ on a finite state space $\mathbb S$ with row-stochastic transition probability matrix $\mathbf P$. It is assumed that the Markov chain is both ergodic (irreducible) and mixing (aperiodic) and homogeneous in time. It follows that there exists a unique stationary distribution π , which satisfies $\pi \mathbf P = \pi$. We take $X_0 = \pi$ so that our Markov chain is stationary, that is, we have $X_n = \pi \mathbf P^n = \pi$ for all $n \in \mathbb Z$. We define the first passage time to a set $\mathbb D \subset \mathbb S$ as

$$\tau_{\mathbb{D}}^+(n) := \inf\{k \ge 0 : X_{n+k} \in \mathbb{D}\},\tag{4}$$

and the last exit time from $\mathbb D$ as

$$\tau_{\mathbb{D}}^{-}(n) := \inf\{k \ge 0 : X_{n-k} \in \mathbb{D}\}. \tag{5}$$

The last exit time is a stopping time with respect to the time-reversed process $(X_{-n})_{n\in \mathbb{Z}}$, namely, the Markov chain on \mathbb{S} with transition probability matrix $\mathbf{P}^- = (P_{ii}^-)$ whose entries are given by

$$P_{ij}^- := \frac{\pi_j}{\pi_i} P_{ji}. \tag{6}$$

Let \mathbb{A} , \mathbb{B} be two nonintersecting subsets of \mathbb{S} such that neither is reachable in one step starting from the other. Following the nomenclature used in physical chemistry literature, at time n we say that the process is forward-reactive $R^+(n)$ [respectively, backward-reactive, $R^-(n)$] according to the realization of the events,

$$R^{\pm}(n) := \left\{ \tau_{\mathbb{B}}^{\pm}(n) < \tau_{\mathbb{A}}^{\pm}(n) \right\}. \tag{7}$$

Then, the process is reactive at time n according to the realization of the event R(n), where

$$R(n) := \{ R^{-}(n) \cup R^{+}(n) \}. \tag{8}$$

In summary, a trajectory is reactive at time n if its most recent visit to $\mathbb{A} \cup \mathbb{B}$ was to \mathbb{A} , it is currently outside of $\mathbb{A} \cup \mathbb{B}$, and its next visit to $\mathbb{A} \cup \mathbb{B}$ will be to \mathbb{B} . One thinks of \mathbb{A} as a source and \mathbb{B} as a target for some process. Associated to the forward and backward reactivities are the forward and backward committors $q_i^{\pm}(n)$ defined for $i \in \mathbb{S}$ by

$$q_i^{\pm}(n) := \Pr\left(R^{\pm}(n) \mid X_n = i\right). \tag{9}$$

We comment briefly on the intuition behind the committors. Note that it is not decidable at time n whether a process is forward reactive at time n. If an ensemble of trajectories each have $X_n = i$, then only some fraction of trajectories will hit $\mathbb B$ before $\mathbb A$; this fraction is exactly $q_i^+(n)$. States with large values of $q_i^+(n)$ tend to be close to $\mathbb B$ in the sense that there is a short path from i to $\mathbb B$ but, in general, this need not be the case. A similar intuition for $q_i^-(n)$ holds

for the time-reversed Markov chain. The committors are the fundamental quantities of transition path theory since they contain all of the information about the (infinite) past and future.

One can show 30 that in the case of a homogeneous and stationary Markov chain, the committors are independent of n and satisfy linear matrix equations

$$q_i^+ = egin{cases} \sum_{j \in \mathbb{S}} P_{ij} q_j^+, & i \notin \mathbb{A} \cup \mathbb{B}, \\ 0, & i \in \mathbb{A}, \\ 1, & i \in \mathbb{B} \end{cases}$$

and

$$q_i^- = \begin{cases} \sum_{j \in \mathbb{S}} P_{ij}^- q_j^-, & i \notin \mathbb{A} \cup \mathbb{B}, \\ 1, & i \in \mathbb{A}, \\ 0, & i \in \mathbb{B}. \end{cases}$$
(10)

Using the committors, a number of statistics can be computed for reactive trajectories. First, we have the reactive density

$$\mu_i^{\mathbb{AB}}(n) := \Pr(X_n = i, R(n)) = q_i^- \pi_i q_i^+, \quad i \notin \mathbb{A} \cup \mathbb{B}. \tag{11}$$

States with large reactive densities relative to their neighbors are interpreted as bottlenecks for reactive trajectories. We also define the reactive current

$$f_{ij}^{\mathbb{AB}}(n) := \Pr(X_n = i, R^-(n), X_{n+1} = j, R^+(n+1)) = q_i^- \pi_i P_{ii} q_i^+, \quad i, j \in \mathbb{S},$$
(12)

as well as the effective reactive current

$$f_{ij}^{+} := \max \left\{ f_{ij}^{\mathbb{AB}} - f_{ji}^{\mathbb{AB}}, 0 \right\}. \tag{13}$$

The effective reactive current is a \mathbb{B} -facing gradient of the reactive density; it identifies pairs of states with a large net flow of probability. Finally, and of particular importance, here is the transition time $t^{\mathbb{A}\mathbb{B}}$. The original definition by Vanden-Eijnden⁸ of $t^{\mathbb{A}\mathbb{B}}$ is as the limiting ratio of the time spent during reactive transitions from \mathbb{A} to \mathbb{B} to the rate of reactive transitions leaving \mathbb{A} . In Ref. 30, the following expression is provided in the discrete case:

$$t^{\mathbb{AB}} := \frac{\Pr(R(n))}{\Pr(R^+(n+1), X_n \in A)} = \frac{\sum_{i \in \mathbb{S}} \mu_i^{\mathbb{AB}}}{\sum_{i \in \mathbb{A}} f_{ii}^{\mathbb{AB}}}.$$
 (14)

In Sec. IV, we will introduce a generalization of the transition time, which will allow us to express Eq. (14) as a straightforward expectation.

C. Open systems and connectivity

In many cases, the trajectory data are given on an open domain, and further processing is required to obtain a suitable Markov chain. We follow Miron *et al.*²⁴ and subsequent works and introduce a two-way nirvana state to create a closed system. Suppose that all trajectory data are contained inside a domain $\mathbb{Y} \subset \mathbb{X}$. We partition \mathbb{Y} as $\mathbb{Y} = \mathbb{Y}^O \cup \omega$ such that $\partial \mathbb{Y} \subset \omega$ and $|\omega| \ll |\mathbb{Y}^O|$. We then construct a covering by N boxes of \mathbb{Y}^O with one additional box appended corresponding to the whole of ω , the nirvana state. Applying Eq. (3), we

obtain a row-stochastic transition matrix of the form

$$\mathbf{P} = \begin{pmatrix} \mathbf{P}^{O \to O} & \mathbf{P}^{O \to \omega} \\ \mathbf{P}^{\omega \to O} & 0 \end{pmatrix},\tag{15}$$

where $\mathbf{P}^{O \to O}$ is $N \times N$, $\mathbf{P}^{O \to \omega}$ is $N \times 1$, and $\mathbf{P}^{\omega \to O}$ is $1 \times N$. Note that trajectories which begin and end in the nirvana state are ignored. In general, we are only interested in reactive trajectories that do not visit this extra nirvana state. This requirement is equivalent to making the replacements $\mathbb{A} \to \mathbb{A} \cup \omega$ and $\mathbb{B} \to \mathbb{B} \cup \omega$ in the basic TPT formulas. One can show²⁴ that this is also equivalent to leaving \mathbb{A} and \mathbb{B} unchanged, but replacing \mathbf{P} with the row-substochastic matrix $\mathbf{P}^{O \to O}$ and π by restriction of the stationary distribution of Eq. (15) to O. We apply the latter method due to the convenience of computation.

Depending on the shape of the data, $\mathbf{P}^{O \to O}$ may not an irreducible, aperiodic matrix. To remedy this, we apply Tarjan's algorithm³¹ to extract the largest strongly connected component of $\mathbf{P}^{O \to O}$. Then, $\mathbf{P}^{O \to O}$ is modified to remove all other states including contributions from trajectories, which pass through the removed states. The final result is that we have an irreducible, aperiodic matrix that avoids the nirvana state and is suitable for use in the formulation of TPT above.

III. TRANSITION TIME STABILITY AND COVERINGS

A. Regular coverings

We apply the methods described in Sec. II to drifter trajectory data obtained from the NOAA Global Drifter Program (GDP).6 In particular, we use quarter-daily interpolated data of the positions of drifters in the tropical Atlantic. We are primarily interested in undrogued drifter trajectories as these may be more accurate models for the motion of Sargassum than their drogued counterparts, as noted in the Introduction. After discarding sections of trajectories which still have their drogue, we must choose a time step T, cf. Sec. II A. Following, e.g., Beron-Vera et al., 26 we choose T = 5 days, a timescale much longer than the Lagrangian decorrelation time scale for the ocean of 1 day.³² This ensures that the assumption of Markovianity will hold to suitable accuracy. In general, Eq. (3) is used except where trajectories contain holes or the length of the trajectory is shorter than T. After obtaining the transition matrix, we apply Eq. (11) to calculate the reactive density, choosing A concentrated off the coast of West Africa (a single box centered at 17° N, 18° W) and B as the Gulf of Mexico (boxes west of 90° W and in [10°N, 30°N]). We cover the computational domain with 760 boxes, resulting in boxes with side lengths of about 2.4°, of which 463 both contained data and were not disconnected. We will sometimes refer to this partition loosely as a partition into "squares," keeping in mind that the covering actually exists on a 2-sphere. Figure 1 shows the distribution of raw counts in the squares. The calculation of the reactive density was repeated with a covering of 780 boxes. The results are shown in Fig. 2.

In addition to the dramatic difference in the density of $\mu^{\mathbb{A}\mathbb{B}}$, we find that $t^{\mathbb{A}\mathbb{B}}=9.81\,\mathrm{yr}$ for the coarser partition (Fig. 2, top panel) and $t^{\mathbb{A}\mathbb{B}}=186\,\mathrm{yr}$ for the finer partition (Fig. 2, bottom panel). In general, changing the number of boxes in the covering results in $\mu^{\mathbb{A}\mathbb{B}}$ graphs that oscillate between patterns similar to the distributions in Fig. 2. We will first address the question of why a small

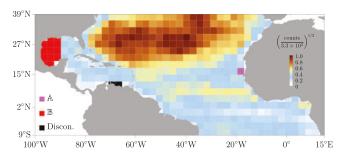
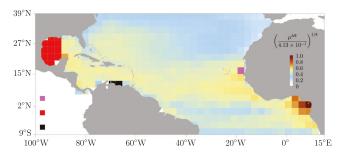


FIG. 1. The number of $x_0(t)$ points falling in each box of the covering. There were 382 793 total trajectories.

change in the number of boxes in the covering can lead to a large change in these TPT statistics. Note that these large changes are not caused by our choice of $\mathbb A$ or $\mathbb B$, it is an issue caused by the nature of coverings by a regular grid of squares. A fundamental issue with a covering by squares is that an addition of even a small number of boxes can result in a shift in the location of every box in the covering. When dealing with sparse trajectory data, this can radically change the outflow in certain regions of the space. This is observed in Fig. 2 in the region $[75^{\circ}\text{W}, 20^{\circ}\text{W}] \times [20^{\circ}\text{N}, 39^{\circ}\text{N}]$. Here, the bottom panel of Fig. 2 shows a much larger portion of reactive density, apparently



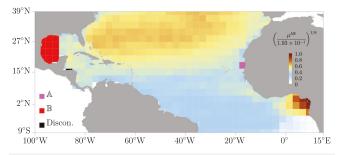


FIG. 2. (Top panel) The eighth-root transformation of the reactive density $\mu^{\mathbb{AB}}$ in the North Atlantic constructed from GDP undrogued drifter data. The computational domain was initialized with 760 square boxes. Boxes colored in black contained data but were removed due to being part of a reducible subset of the directed graph associated with the Markov chain resulting from discretizing the drifter motion using Ulam's method. (Bottom panel) As in the top panel, but with an initialization of 780 boxes.

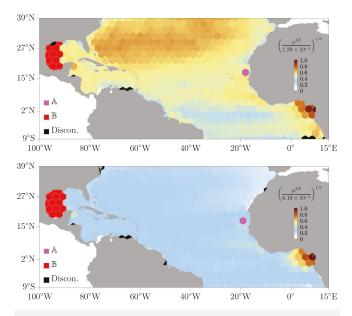


FIG. 3. As in Fig. 2 but with a hexagonal covering. The computational domain in the top (respectively, bottom) panel was initialized with 820 (respectively, 840) boxes.

suggesting that particles tend to circulate in this area before eventually finding the more direct path from $\mathbb A$ to $\mathbb B$ highlighted in the top panel of Fig. 2. Note also that this region is sparsely populated as shown in Fig. 1.

Another option for a regular covering is by hexagons. Hexagons could be considered a more natural choice than squares since the distance between the centers of adjacent hexagons is constant. In Ref. 33, a hexagonal covering provided by the H3 spatial index34 was used to construct the transition matrix. We choose the same parameters as for the squares, but instead cover the computational domain by a regular grid of hexagons. This is repeated twice with a small difference in the number of initial cells; the results are shown in Fig. 3. Again, we find that a small change in the number of covering cells leads to a large change in transition path theory statistics. In addition to the differences in the reactive densities, we find $t^{AB} = 25.6 \,\text{yr}$ for the coarser partition (Fig. 3, top panel) and $t^{\mathbb{AB}} = 156$ yr for the finer partition (Fig. 3, bottom panel). The essential problem is the same for both square and hexagonal coverings, namely, that it is not clear which resulting statistics should be trusted. One would hope that variations in a scalar statistic such as t^{AB} would settle as the number of boxes increases, but this is not the case with regular coverings. Motivated by these examples, we propose another kind of covering that addresses this issue.

B. Voronoi coverings

We propose that instead of covering the domain with a regular grid, we instead cluster the observations and draw polygons based on the cluster boundaries. The intended result should be that data points that are close to each other should tend to end up in the same box, and, hence, the derived transition matrix should be robust

against small changes in box number. There are a number of clustering algorithms; for the current application, we have chosen to use the k-means method. 35,36 This is a hard clustering algorithm, which is guaranteed to converge and create n clusters (if possible) when requested. In addition, k-means is straightforward to implement and is built into the clustering packages of many popular languages. The output of k-means is a collection of centroids such that each data point belongs to the cluster defined by its closest centroid in terms of Euclidean distance. Hence, we can define the polygons covering the computational domain using a Voronoi tessellation; cf. e.g., Burrough, McDonnell, and Lloyd.³⁷ While this kind of technique is used in the chemical literature, 38 we emphasize that our goal in generating these clusters is to increase the stability of TPT statistics. We compute the intersection of the convex hull of the data with the Voronoi tessellation to reduce the size of the outer cells for clearer visualization. Performing the same reactive density calculation as in Sec. III A gives the results shown in Fig. 4. The picture is rather insensitive to the number of boxes considered, which we quantify below.

For this covering, we find $t^{\mathbb{A}\mathbb{B}}=2.34\,\mathrm{yr}$. Note that Fig. 4 does not contain any disconnected polygons, and, in general, the Voronoi covering is less prone to disconnections although they are still possible. To compare the stability of this method to the regular coverings discussed previously, we compute $t^{\mathbb{A}\mathbb{B}}$ for a number of boxes sizes between 20 and 600 as shown in Fig. 5. We see that the Voronoi covering produces significantly more stable transition times of consistently reasonable magnitudes. To understand this, we note that the addition of a small number of clusters does not tend to have a large global effect as observed for regular coverings. Requesting more clusters tends to subdivide larger clusters or create more where the data are dense and leave others untouched.

The Voronoi covering has some drawbacks. First, although the k-means algorithm is a relatively fast clustering algorithm and amenable to parallelization if necessary, it is still roughly two orders of magnitude slower computationally than the regular coverings. This first attempt may be improved by alternate clustering algorithms; for instance, Prinz *et al.*³⁹ studied faster clustering algorithms with comparable results. In addition, the initial guess for the location of centroids in the typical k-means algorithm is random. We find that the variation in $t^{\mathbb{AB}}$ caused by variations in this initial guess is roughly five steps on average, much smaller than those coming from the change in box size. We recommend running the clustering

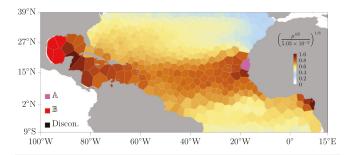


FIG. 4. As in Fig. 2, but with a Voronoi covering generated by k-means with 500 clusters

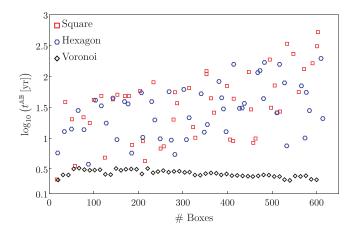
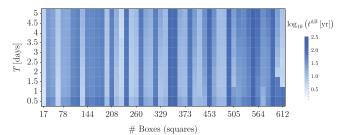


FIG. 5. The transition time of Eq. (14) for various box sizes and polygon coverings. The horizontal axis shows the number of boxes remaining after boxes with no data were removed, e.g., 600 boxes for a square covering is the result of an initial covering of 1000 boxes.

algorithm multiple times to check that the initial guess has not accidentally found an undesirable local maximum. Finally, we note that the nature of the algorithm means that it is difficult to add boxes in specific locations, but this also applies for non-adaptive regular coverings.

C. The time step T

As mentioned in Sec. III A, the time step T = 5 days was chosen based on time scales arising from the Lagrangian characteristics (decorrelation) in the upper ocean. Here, we explore the validity of this choice based on the stability of the transition time. Figure 6 shows this transition time as a function of both box number and time step T both for a regular square covering and the Voronoi covering. Reading across the horizontal axis for a fixed T, we see the same results as in Fig. 5, namely, that the transition time is stable against box number for the Voronoi covering but not for the regular covering. For a fixed box number, reading up the vertical axis generally shows that T = 0.5 days tends to have a slightly higher t^{AB} but for T > 0.5 days, there is very little variation for both the regular and Voronoi coverings. When T = 0.5 days, there are enough trajectories that do not leave their initial cells that the transition matrix is very strongly diagonal; this serves to increase the transition time. As discussed previously, small changes in the number of boxes for a regular covering can result in global shifts in the locations of boxes in the covering. However, small changes in T do not generally have this behavior since increasing T still leaves the same number of observations (modulo a small number of points left off at the end) and, hence, for sufficiently long trajectories, the effect is not felt to a significant degree. With real data, there will, of course, be an upper limit to T beyond which results cease to become trustworthy due to lack of communication between boxes and large numbers of short trajectories being rejected. What we observe here, regardless of considerations related to the Lagrangian decorrelation time of the



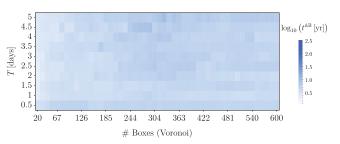


FIG. 6. (Top panel) The transition time of Eq. (14) with a regular covering of squares for various box sizes and time steps. The horizontal axis shows the number of boxes after boxes with no data removed. For a regular hexagonal covering, the graph looks very similar. (Bottom panel) As in the top panel, but using the Voronoi covering described in Sec. III B.

ocean, is that the lower limit for obtaining stable transition times is roughly 1 day.

IV. A GENERALIZED TRANSITION TIME

In this section, we present a generalization of Eq. (14), which can be used to obtain a partition of the computational domain based on the time it takes to reach $\mathbb B$ from an arbitrary cell. The aim of this generalization is to provide local information, that is, information about reactive trajectories at a particular state that have already left the source $\mathbb A$. Let

$$\mathbb{C}^{+} = \left\{ i \notin \mathbb{B} : \sum_{\ell \in \mathbb{S}} P_{i\ell} q_{\ell}^{+} > 0 \right\}. \tag{16}$$

We define the *remaining time* $t^{\mathbb{B}}$ for all $n \in \mathbb{Z}$ as

$$t^{i\mathbb{B}} := \begin{cases} \operatorname{Ex}\left[\tau_{\mathbb{B}}^{+}(n+1) \mid X_{n} = i, \ R^{+}(n+1)\right], & i \notin \mathbb{B}, \\ 0, & i \in \mathbb{B}. \end{cases}$$
(17)

A similar formula is referred to as the *lead time* in Finkel *et al.*⁴⁰ In the Appendix, we establish the following Lemma:

Lemma 1: Equation (17) satisfies a set of linear equations,

$$t^{i\mathbb{B}} := \begin{cases} 1 + \sum_{j \in \mathbb{C}^+} \frac{P_{ij}q_j^+}{\sum_{\ell \in \mathbb{S}} P_{i\ell}q_\ell^+} t^{jB}, & i \in \mathbb{C}^+, \\ 0, & i \in \mathbb{B}. \end{cases}$$
(18)

When A contains only one state, we also have that

$$t^{\mathbb{IB}}\big|_{i=\mathbb{A}} = t^{\mathbb{AB}} + 1,\tag{19}$$

where $t^{\mathbb{AB}}$ is defined in Eq. (14).

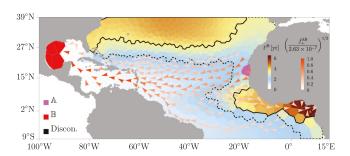


FIG. 7. The remaining time of Eq. (17) for a Voronoi covering with 500 cells with the effective reactive current of Eq. (13) overlaid. The remaining time can be partitioned into three regions via k-means clustering such that states from each region have similar remaining times. The regions are shown as level sets demarcated by the solid and dashed black lines. The average remaining times in each region are: 4.0 yr (solid line interior), 2.8 yr (between solid and dashed lines), and 1.3 yr (outside both solid and dashed lines).

By applying Lemma 1, we can compute the remaining time for each box in a given covering. We choose a Voronoi covering with 500 boxes, similar to the construction of Fig. 4. To build a remainingtime-based dynamical geography, we partition the remaining times into three clusters via k-means. We show the outline of this geography overlaid with the effective reactive current of Eq. (13) in Fig. 7. The dynamical geography obtained here is similar to the one obtained in Beron-Vera et al.7 by other means. We see that the longest times are found near the Gulf of Guinea and the most subtropical North Atlantic. This is consistent with the drifter data: there is a large inflow to the Gulf of Guinea, making drifters near the West coast of Africa cause a large pileup of trajectories in this region. Similarly, the Gulf Stream pushes drifters up and out of the Gulf of Mexico such that they are unlikely to transition back into the Gulf in a short time once they pass the coast of Florida. Consequently, the dynamical geography provided by the remaining time is related to the distance between states and the target, but they are not interchangeable; the remaining time provides additional information.

Examining the reactive current, we recover the two main transition paths observed in Beron-Vera *et al.*, namely, the direct westward path and the indirect path which initially moves toward the Gulf of Guinea before circulating westward across the equatorial Atlantic. Not only does the direct path have a larger effective current, the transition times are also shorter. In general, there is a noticeable separation between westward and eastward-bound currents south of the source A. We note here that the remaining time of a state need not be positively correlated with its effective current. For example, the effective currents are roughly equal in the Gulf of Guinea and northern portion of the Gulf of Mexico but the remaining times are significantly different.

V. CONCLUSIONS

When Ulam's method is applied to trajectory data, the space must be partitioned into a covering to discretize the motion and thereby construct a transition probability matrix. We have shown that two types of standard coverings made up of regular grids of

squares and hexagons result in unstable transition times when Transition Path Theory (TPT) is applied to the induced Markov chain. Changing the number of squares or hexagons in the covering leads to global shifts in their location, producing untrustworthy results for TPT statistics. We proposed a different kind of covering that partitions the space into Voronoi cells based on k-means clustering of the observations. This covering leads to transition times, which are stable against the number of requested clusters. This algorithm was chosen for simplicity and effectiveness, but there are many clustering algorithms that could be explored, in particular, with consideration toward improving computational performance for large data sets. In addition, we found that the transition time does not depend strongly on the time step through which the trajectory data are sliced for time steps between 1 and 5 days for undrogued drifters in the tropical/subtropical Atlantic. Finally, we introduced a generalization of the standard TPT transition time, which contains the standard TPT transition time as a special case. Clustering cells based on this generalized transition time produce a partition of the domain, which reveals weakly dynamically connected regions.

ACKNOWLEDGMENTS

The authors are grateful to Luzie Helfmann for providing notes which inspired the development of Eq. (17). This work was supported by the National Science Foundation (NSF) under Grant No. OCE2148499.

AUTHOR DECLARATIONS

Conflict of Interest

The authors have no conflicts to disclose.

Author Contributions

G. Bonner: Conceptualization (equal); Formal analysis (lead); Software (lead); Visualization (lead); Writing – original draft (lead); Writing – review & editing (equal). **F. J. Beron-Vera:** Conceptualization (equal); Funding acquisition (equal); Writing – review & editing (equal). **M. J. Olascoaga:** Conceptualization (equal); Funding acquisition (equal); Writing – review & editing (equal).

APPENDIX: PROOF OF LEMMA 1

We first establish that Eq. (17) can be written as the solution to the system of linear equations in Eq. (18). In what follows, we repeatedly use the Markov property and the stationarity of our chain. First, we have that

$$Pr(X_{n+1} = j \mid X_n = i, R^+(n+1))$$

$$= \frac{Pr(X_{n+1} = j, R^+(n+1) \mid X_n = i)}{Pr(R^+(n+1) \mid X_n = i)}$$
(A1)

$$=\frac{P_{ij}q_j^+}{\sum_{\ell\in\mathbb{S}}P_{i\ell}q_\ell^+}.$$
 (A2)

Note that the condition that our Markov chain is ergodic does not necessarily imply that $\sum_{\ell \in S} P_{i\ell} q_\ell^+ \neq 0$ for all $i \notin \mathbb{B}$. There can exist a series of states for which the only path between them and \mathbb{B} passes

through \mathbb{A} . Hence, an "interior" state whose neighbors all have $q_\ell^+=0$ will have $\sum_{\ell\in S}P_{i\ell}q_\ell^+\neq 0$. To address this case, we introduce the set of states

$$\mathbb{C}^{+} = \left\{ i \notin \mathbb{B} : \sum_{\ell \in \mathbb{S}} P_{i\ell} q_{\ell}^{+} > 0 \right\}$$
 (A3)

and restrict $t^{\mathbb{I}\mathbb{B}}$ to these states. Therefore, Eq. (A2) is well-defined. Taking $i \notin \mathbb{B}$ in Eq. (17), we condition on the value of X_{n+1} and use the fact that $\tau_B^+(n+1) = 1 + \tau_B^+(n+2)$ on the event that $X_{n+1} \notin \mathbb{B}$ to obtain

$$t^{\mathbb{IB}} := \begin{cases} 1 + \sum_{j \in \mathbb{C}^+} \frac{P_{ij}q_j^+}{\sum_{\ell \in \mathbb{S}} P_{i\ell}q_\ell^+} t^{jB}, & i \in \mathbb{C}^+, \\ 0, & i \in \mathbb{B}. \end{cases}$$
(A4)

We will now show that Eq. (17) with $i = \mathbb{A}$ coincides with Eq. (14). We compute the quantity $Z(n) = \Pr(R^-(n), R^+(n))$. For the time n to be reactive, there must be a last visit to \mathbb{A} at some time $\ell < n$, and the next visit to $\mathbb{A} \cup \mathbb{B}$ must be to \mathbb{B} at least $n - \ell$ steps later. Therefore, we can write

$$Z(n) = \sum_{\ell < n} \Pr\left(X_{\ell} \in \mathbb{A}, R^{+}(\ell+1), \tau_{\mathbb{B}}^{+}(\ell+1) > n-\ell\right)$$

$$= \sum_{\ell < n} \Pr\left(\tau_{\mathbb{B}}^{+}(\ell+1) \ge n-\ell \mid X_{\ell} \in \mathbb{A}, R^{+}(\ell+1)\right)$$
(A)

$$\times \Pr\left(X_{\ell} \in \mathbb{A}, R^{+}(\ell+1)\right). \tag{A6}$$

By the stationarity of our process, $\Pr(X_{\ell} \in \mathbb{A}, R^{+}(\ell + 1))$ is independent of ℓ , so we have

$$Z(n) = \Pr\left(X_n \in \mathbb{A}, R^+(n+1)\right)$$

$$\times \sum_{\ell \in n} \Pr\left(\tau_{\mathbb{B}}^+(\ell+1) > n-\ell \mid X_\ell \in \mathbb{A}, R^+(\ell+1)\right). \tag{A7}$$

Applying the stationary property again and re-indexing the sum over ℓ gives

$$Z(n) = \Pr\left(X_n \in \mathbb{A}, R^+(n+1)\right)$$

$$\times \sum_{\ell=1}^{\infty} \Pr(\tau_{\mathbb{B}}^+(n+1) > \ell \mid X_n \in \mathbb{A}, R^+(n+1)). \tag{A8}$$

Since for a nonnegative discrete random variable *X* we have $\text{Ex}[X] = \sum_{k>1} P(X \ge k)$, we conclude that

$$Z(n) = \Pr\left(X_n \in \mathbb{A}, R^+(n+1)\right)$$

$$\times \left(-1 + \operatorname{Ex}\left[\tau_{\mathbb{R}}^+(n+1) \mid X_n = i, R^+(n+1)\right]\right), \quad (A9)$$

which implies that Vanden-Eijden's⁴¹ $t^{\mathbb{A}\mathbb{B}}$ and Eq. (17) with $i = \mathbb{A}$ differ by one step, i.e., they are identical except that $t^{\mathbb{A}\mathbb{B}}$ does not "count" the first step to leave \mathbb{A} .

DATA AVAILABILITY

The data employed in this paper are openly available from the NOAA Global Drifter Program at http://www.aoml.noaa.gov/phod/dac/. The computations were carried out using Julia; a package has been developed, which is distributed from https://github.com/70Gage70/UlamMethod.jl.

REFERENCES

- ¹J. N. Butler, B. F. Morris, J. Cadwallader, and A. W. Stoner, "Studies of Sargassum and the Sargassum community," Bermuda Biol. Station Spec. Publ. **22**, 307 (1983)
- ²J. J. Milledge and P. J. Harvey, "Golden tides: Problem or golden opportunity? The valorisation of *Sargassum* from beach inundations," J. Mar. Sci. Eng. **4**, 60 (2016).
- ³D. D. Laffoley, H. S. Roe, M. Angel, J. Ardron, N. Bates, I. Boyd, S. Brooke, K. N. Buck, C. Carlson, B. Causey *et al.*, "The protection and management of the Sargasso Sea," Technical Report (Sargasso Sea Alliance, 2011).
- ⁴J. Gower, E. Young, and S. King, "Satellite images suggest a new *Sargassum* source region in 2011," Remote Sens. Lett. **4**, 764–773 (2013).
- ⁵B. van Tussenbroek, H. Arana, R. Rodriguez-Martinez, J. Espinoza-Avalos, H. Canizales-Flores, C. Gonzalez-Godoy, M. Barba-Santos, A. Vega-Zepeda, and L. Collado-Vides, "Severe impacts of brown tides caused by *Sargassum* spp. on near-shore Caribbean seagrass communities," Mar. Pollut. Bull. **122**, 272–281 (2017).
- ⁶R. Lumpkin and M. Pazos, "Measuring surface currents with surface velocity program drifters: The instrument, its data and some recent results," in *Lagrangian Analysis and Prediction of Coastal and Ocean Dynamics*, edited by A. Griffa, A. D. Kirwan, A. Mariano, T. Özgökmen, and T. Rossby (Cambridge University Press, 2007), Chap. 2, pp. 39–67.
- ⁷F. J. Beron-Vera, M. J. Olascoaga, N. F. Putman, J. Trinanes, R. Lumpkin, and G. Goni, "Dynamical geography and transition paths of *Sargassum* in the tropical Atlantic," AIP Adv. **12**, 105107 (2022).
- ⁸E. Vanden-Eijnden, "Transition path theory," in *Computer Simulations in Condensed Matter Systems: From Materials to Chemical Biology* (Springer, 2006), Vol. 1, pp. 453–493.
- ⁹W. E and E. Vanden-Eijnden, "Transition-path theory and path-finding algorithms for the study of rare events," Annu. Rev. Phys. Chem. 61, 391–420 (2010).
 ¹⁰M. Wang, C. Hu, B. Barnes, G. Mitchum, B. Lapointe, and J. P. Montoya, "The great Atlantic Sargassum belt," Science 365, 83–87 (2019).
 ¹¹J. Franks, D. Johnson, and D. Ko, "Pelagic Sargassum in the tropical North
- ¹¹J. Franks, D. Johnson, and D. Ko, "Pelagic *Sargassum* in the tropical North Atlantic," Gulf Caribb. Res. **27**, C6–11 (2016).
- ¹² A. Lasota and M. C. Mackey, Chaos, Fractals and Noise: Stochastic Aspects of Dynamics, 2nd ed., Applied Mathematical Sciences Vol. 97 (Springer, New York, 1994).
- ¹³S. M. Ulam, A Collection of Mathematical Problems (Interscience Publishers, 1960), Vol. 8.
- 14T.-Y. Li, "Finite approximation for the Frobenius-Perron operator. A solution to Ulam's conjecture," J. Approx. Theory 17, 177–186 (1976).
 15 J. N. Reddy, Introduction to the Finite Element Method (McGraw-Hill Educa-
- ¹³J. N. Reddy, Introduction to the Finite Element Method (McGraw-Hill Education, 2019).
- ¹⁶P. Miron, F. J. Beron-Vera, M. J. Olascoaga, and P. Koltai, "Markov-chain-inspired search for MH370," Chaos 29, 041105 (2019).
- ¹⁷S. Klus, P. Koltai, and C. Schütte, "On the numerical approximation of the Perron-Frobenius and Koopman operator," J. Comput. Dyn. 3, 51–79 (2016).
- ¹⁸G. Froyland, G. A. Gottwald, and A. Hammerlindl, "A computational method to extract macroscopic variables and their dynamics in multiscale systems," SIAM J. Appl. Dyn. Syst. **13**, 1816–1846 (2014).
- ¹⁹E. Van Sebille, M. H. England, and G. Froyland, "Origin, dynamics and evolution of ocean garbage patches from observed surface drifters," Environ. Res. Lett. 7, 044040 (2012).
 ²⁰O. Junge and P. Koltai, "Discretization of the Frobenius-Perron operator using
- ²⁰O. Junge and P. Koltai, "Discretization of the Frobenius-Perron operator using a sparse HAAR tensor basis: The sparse Ulam method," SIAM J. Numer. Anal. 47, 3464–3485 (2009).
- ²¹P. Miron, F. J. Beron-Vera, M. J. Olascoaga, J. Sheinbaum, P. Pérez-Brunius, and G. Froyland, "Lagrangian dynamical geography of the Gulf of Mexico," Sci. Rep. 7, 7021 (2017).
- ²²M. J. Olascoaga, P. Miron, C. Paris, P. Pérez-Brunius, R. Pérez-Portela, R. H. Smith, and A. Vaz, "Connectivity of Pulley Ridge with remote locations as inferred from satellite-tracked drifter trajectories," J. Geophys. Res. 123, 5742–5750, https://doi.org/10.1029/2018JC014057 (2018).
- ²³ P. Miron, F. J. Beron-Vera, M. J. Olascoaga, G. Froyland, P. Pérez-Brunius, and J. Sheinbaum, "Lagrangian geography of the deep Gulf of Mexico," J. Phys. Oceanogr. 49, 269–290 (2019).

- ²⁴P. Miron, F. Beron-Vera, L. Helfmann, and P. Koltai, "Transition paths of marine debris and the stability of the garbage patches," Chaos 31, 033101
- (2021). ²⁵F. J. Beron-Vera, N. Bodnariuk, M. Saraceno, M. J. Olascoaga, and C. Simionato, "Stability of the Malvinas current," Chaos 30, 013152 (2020).
- ²⁶F. J. Šeron-Vera, M. J. Olascoaga, N. F. Putman, J. Triñanes, R. Lumpkin, and G. Goni, "Dynamical geography and transition paths of Sargassum in the tropical Atlantic," AIP Adv. 105107, 105107 (2022).
- ²⁷M. J. Olascoaga and F. J. Beron-Vera, "Exploring the use of Transition Path Theory in building an oil spill prediction scheme," Frontiers 9, 1041005 (2023).
- ²⁸F. J. Beron-Vera, M. J. Olascoaga, L. Helfmann, and P. Miron, "Samplingdependent transition paths of Iceland-Scotland overflow water," J. Phys. ceanogr. 53, 1151-1160 (2023).
- ²⁹P. Metzner, C. Schütte, and E. Vanden-Eijnden, "Illustration of transition path theory on a collection of simple examples," J. Chem. Phys. 125, 084110
- ³⁰L. Helfmann, E. Ribera Borrell, C. Schütte, and P. Koltai, "Extending transition path theory: Periodically driven and finite-time dynamics," J. Nonlinear Sci. 30, 3321-3366 (2020).
- ³¹R. Tarjan, "Depth-first search and linear graph algorithms," SIAM J. Comput. 1, 146-160 (1972).
- ³²J. LaCasce, "Statistics from Lagrangian observations," Prog. Oceanogr. 77, 1–29 (2008).

- 33 M. O'Malley, A. M. Sykulski, R. Laso-Jadart, and M.-A. Madoui, "Estimating the travel time and the most likely path from Lagrangian drifters," J. Atmos. Oceanic Technol. 38, 1059-1073 (2021).
- ³⁴UBER, "H3 spatial index" (2019), see https://eng.uber.com/h3/ (last accessed December 1, 2022).
- 35 E. W. Forgy, "Cluster analysis of multivariate data: Efficiency versus inter-
- pretability of classifications," Biometrics 21, 768–769 (1965).

 36S. Lloyd, "Least squares quantization in PCM," IEEE Trans. Inf. Theory 28, 129-137 (1982).
- ³⁷P. A. Burrough, R. A. McDonnell, and C. D. Lloyd, *Principles of Geographical* Information Systems (Oxford University Press, 2015).
- ³⁸F. Noé, C. Schütte, E. Vanden-Eijnden, L. Reich, and T. R. Weikl, "Constructing the equilibrium ensemble of folding pathways from short off-equilibrium simulations," Proc. Natl. Acad. Sci. U. S. A. 106, 19011-19016 (2009).
- ³⁹J.-H. Prinz, H. Wu, M. Sarich, B. Keller, M. Senne, M. Held, J. D. Chodera, C. Schütte, and F. Noé, "Markov models of molecular kinetics: Generation and validation," J. Chem. Phys. 134, 174105 (2011).
- ⁴⁰ J. Finkel, R. J. Webber, E. P. Gerber, D. S. Abbot, and J. Weare, "Learning forecasts of rare stratospheric transitions from short simulations," Mon. Weather Rev. **149**, 3647-3669 (2021).
- ⁴¹E. Vanden-Eijnden, "Transition path theory," in Computer Simulations in Condensed Matter Systems: From Materials to Chemical Biology, edited by M. Ferrario, G. Ciccotti, and K. Binder (Springer, Berlin, 2006), Vol. 1, pp. 453-493.