# Exploring sudden stratospheric warmings with transition path theory
## Supplementary information

Justin Finkel, Edwin P. Gerber, Dorian S. Abbot, Jonathan Weare

February 4, 2022

This document has three sections. Section 1 spells out transition path theory (TPT) formally, with definitions and equations for all quantities of interest without regard to their numerical approximation. Section 2 describes the numerical method, dynamical Galerkin Approximation (DGA), and provides some numerical benchmarks. Section 3 then gives details about the optimization method we used to find a minimum-action path.

# 1    Transition path theory formalism

## 1.1    The transition path ensemble

We begin a quantitative description of transition paths by formalizing the ensemble notion. The theoretical development parallels [612006Vanden-Eijnden], but expands on it in several ways. Consider the stratosphere, or any other stochastic ergodic dynamical system, evolving through a very long time interval $(-T, T)$, during which it crosses from $A$ to $B$ and back a number $M_T$ of times. As $T \to \infty$, ergodicity guarantees that $M_T \to \infty$ as well. The $m$th transition path begins at time $\tau_m^-$ (so $\mathbf{X}(\tau_m^-) \in A$[1]) and ends at time $\tau_m^+$ (so $\mathbf{X}(\tau_m^+) \in B$). Each $\tau_m^-$ marks the beginning of an orange segment in Fig. 2, and $\tau_m^+$ marks the end of it.

We will describe the transition path ensemble at two levels of granularity. At the first level, we consider the set of *reactive snapshots*, which are the instantaneous model states $\mathbf{X}(t)$ realized in the course of a transition without regard to their ordering in time or their grouping into separate transition events:

$$\text{Reactive snapshots} = \bigcup_{m=1}^{\infty} \bigcup_{t=\tau_m^-}^{\tau_m^+} \mathbf{X}(t). \tag{1}$$

We use "reactive" for consistency with the chemistry literature. This ensemble is described completely by the *reactive densities* such as $\pi_{AB}$ and $\pi_{BA}$ in Fig. 2(b,d), which lump together all reactive snapshots into a single distribution, regardless of which transition event they came from. In other words, every orange point in Fig. 1(c,d) contributes one sample to the ensemble. The reactive densities describe how the stratosphere looks at a single moment in transition, but not how it evolves through the transition.

At the second level, we distinguish each transition path as a unique, coherent object, containing a sequence of snapshots ordered in time. We formally define the $(A \to B)$ transition path ensemble as

$$\text{Transition paths} = \tag{2}$$

$$\left\{ \left\{ (t, \mathbf{X}(t)), \tau_m^- \le t < \tau_m^+ \right\}, m = 1, 2, \dots \right\} \tag{3}$$

The inner set is the collection of snapshots along the $m$th transition path, and the outer set is the collection of paths, which becomes infinite as $T \to \infty$. There is no fixed duration of transition paths; each one has a different duration $\tau_m^+ - \tau_m^-$. For this reason, the space of paths has infinite dimension. Therefore, there is no probability density to describe the path ensemble. However, *functionals* of transition paths do have well-defined distributions. Using the abbreviation $\mathbf{X}^{(m)} := \left\{ (t, \mathbf{X}(t)) : \tau_m^- \le t < \tau_m^+ \right\}$ for the $m$th transition path, we can define arbitrary functionals $\mathscr{G}$ such as

$$\mathscr{G}_1[\mathbf{X}^{(m)}] = \tau_m^+ - \tau_m^-, \tag{4}$$

$$\mathscr{G}_2[\mathbf{X}^{(m)}] = \int_{\tau_m^-}^{\tau_m^+} \overline{v'T'}(30\,\text{km})\, dt, \tag{5}$$

$$\mathscr{G}_3[\mathbf{X}^{(m)}] = \max \left\{ \left| \frac{U(30\,\text{km})(t_2) - U(30\,\text{km})(t_1)}{t_2 - t_1} \right| : \tau_m^- \le t_1 < t_2 \le \tau_m^+ \right\}, \tag{6}$$

---

[1]Technically, we assume $\mathbf{X}(t)$ is right-continuous with left limits, meaning $\mathbf{X}(\tau_m^-) \notin A$ but $\lim_{t \uparrow \tau_m^-} \mathbf{X}(t) \in A$. This detail is not important for us here.

which quantify the elapsed time, the total heat flux at 30 km, and the fastest drop in wind speed at 30 km recorded over the whole transition path. (A time derivative is not technically well-defined for a white noise-driven process.) The quantities of interest $\mathscr{G}$ will, of course, depend on the application. For any fixed scalar-valued or vector-valued $\mathscr{G}$, the collection of random variables $\{\mathscr{G}[\mathbf{X}^{(1)}], \mathscr{G}[\mathbf{X}^{(2)}], \ldots\}$ has a well-defined steady-state distribution that we wish to characterize. In principle, we could do this by "direct numerical simulation" (DNS): integrate the system for a long time, collect many $A \to B$ transition paths $\mathbf{X}^{(m)}$, calculate $\mathscr{G}[\mathbf{X}^{(m)}]$ for each one, and estimate summary statistics. Although DNS is simple and general, it is expensive for high-dimensional models, particularly for rare event simulation. The DGA method, explained below in section 2, circumvents DNS by using only short trajectories (20 days long in our implementation). These are short not only compared to the return time $\tau_{m+1}^- - \tau_m^+$ (~ 1800 days for the Holton-Mass model), but even compared to the ($A \to B$) transit time $\tau_m^+ - \tau_m^-$ (~ 80 days for the Holton-Mass model). We focus on functionals of the *transition path integral* form:

$$\mathscr{G}[\mathbf{X}^{(m)}] = \int_{\tau_m^-}^{\tau_m^+} \Gamma(\mathbf{X}(t)) \, dt \tag{7}$$

where $\Gamma$ is a user-defined quantity of interest. For example $\Gamma(\mathbf{x}) = 1$ yields the transit time (4) and $\Gamma = \overline{v'T'}(30 \text{ km})$ yields the total heat flux at 30 km, (5). For certain extreme weather events, $\Gamma$ might be chosen to measure accumulated damage of some kind, say, the total rainfall deposited over an area (in the case of hurricanes) or total time with surface temperatures above a certain threshold (in the case of heat waves). In a downward-coupled SSW model, one could define $\Gamma$ to reflect the human impact of extreme cold spells, but in the simple Holton-Mass model we value $\Gamma$ only for dynamical insight into SSW variability. We will refer to $\int \Gamma(\mathbf{X}(t)) \, dt$ as a *transition path integral*. A limitation of DGA is that it cannot directly handle more complex nonlinear functionals of the form (6). We present the mathematics of generalized rates in section 2, below. While conceptually promising to explore many aspects of the transition path ensemble, the estimates with our current method are not yet reliable, being "second-order" calculations, which we explain below. Numerical estimation of transition path integrals will be explored in future work.

### 1.1.1 Forecast functions

The essential insight of TPT is to express the quantities of interest in terms of a set of *forecast functions*. A forecast is an estimate of the future conditioned on the present, which in probability language takes the form

$$F^+(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[Q(\{(t, \mathbf{X}(t)) : t \geq 0\})]. \tag{8}$$

Here, $\mathbb{E}_{\mathbf{x}}$ indicates a conditional expectation given a fixed initial condition $\mathbf{X}(0) = \mathbf{x}$ (we can set $t_0 = 0$ when assuming autonomous dynamics). $Q$ is a generic functional of the future evolution of the state $\mathbf{X}(t)$. It is explicitly a random variable under the stochastic forcing we impose here, but even in a deterministic model, uncertainty from initial conditions and model error lead to effective randomness. For example, $Q$ could return 1 if $\mathbf{X}(t)$ next hits $B$ before $A$, and 0 if $\mathbf{X}(t)$ next hits $A$ before $B$. This makes $F^+$ simply the forward committor, as introduced in section b:

$$F^+(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}\left[\mathbb{1}_B(\mathbf{X}(\tau_{A \cup B}^+))\right] \tag{9}$$

$$= \mathbb{P}_{\mathbf{x}}\{\mathbf{X}(\tau_{A \cup B}^+) \in B\} =: q_B^+(\mathbf{x}) \tag{10}$$

We might also wish to forecast the time it takes to get there, by defining $Q = \tau_{A \cup B}^+ \mathbb{1}_B(\mathbf{X}(\tau_{A \cup B}^+))$, which then gives us the lead time $\eta_B^+(\mathbf{x}) = F^+(\mathbf{x})/q_B^+(\mathbf{x})$.

As explained in section b, the forward committor only looks to the future, and the backward committor is needed to distinguish the $A \to B$ phase from the $B \to B$ phase.

$$q_A^-(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[\mathbb{1}_A(\mathbf{X}(\tau_{A \cup B}^-))] = \mathbb{P}_{\mathbf{x}}\{\mathbf{X}(\tau_{A \cup B}^-) \in A\} \tag{11}$$

This is a backward-in-time forecast, or *aftcast*.

Forward and backward committors are central components in the existing transition path theory laid out in [162006E and Vanden-EijndenE, 612006Vanden-Eijnden, 402009Metzner et al.Metzner, Schutte,, and Vanden-Eijnden, 392006Metzner et al.Metzner, Schutte,, and Vanden-E ?], and elsewhere. Here, we generalize committors to forecast not only where the trajectory ends up, but what happens along the way. We consider forecast/aftcast functions of the form

$$F_\Gamma^+(\mathbf{x}; \lambda) = \mathbb{E}_{\mathbf{x}}\left[\mathbb{1}_B(\mathbf{X}(\tau_{A \cup B}^+)) \exp\left(\lambda \int_0^{\tau_{A \cup B}^+} \Gamma(\mathbf{X}(r)) \, dr\right)\right] \tag{12}$$

$$F_\Gamma^-(\mathbf{x}; \lambda) = \mathbb{E}_{\mathbf{x}}\left[\mathbb{1}_A(\mathbf{X}(\tau_{A \cup B}^-)) \exp\left(\lambda \int_{\tau_{A \cup B}^-}^0 \Gamma(\mathbf{X}(r)) \, dr\right)\right] \tag{13}$$

2

where $\lambda$ is a real free parameter. With $\lambda = 0$ (and any $\Gamma$), $F_\Gamma^+$ and $F_\Gamma^-$ in (12) and (13) reduce to the forward and backward committors. With $\Gamma \equiv 1$, $F_\Gamma^+$ and $F_\Gamma^-$ are moment-generating functions for the next hitting time $\tau_{A\cup B}^+$ and the most recent hitting time $\tau_{A\cup B}^-$. We can use the fact that the past and future are independent *conditional on the present*, and multiply them together to get an expectation over transition paths crossing through $\mathbf{x}$:

$$F_\Gamma^-(\mathbf{x}, \lambda) F_\Gamma^+(\mathbf{x}, \lambda) = \tag{14}$$

$$\mathbb{E}_\mathbf{x}\left[ \mathbb{1}_A(\mathbf{X}(\tau_{A\cup B}^-)) \mathbb{1}_B(\mathbf{X}(\tau_{A\cup B}^+)) \exp\left( \lambda \int_{\tau_{A\cup B}^-}^{\tau_{A\cup B}^+} \Gamma(\mathbf{X}(r))\, dr \right) \right]$$

This product is now seen to be a moment-generating function for the transition path integral (7). Differentiating in $\lambda$ at $\lambda = 0$ provides us with any moment of the path integral distribution:

$$\partial_\lambda^k \left[ F^-(\mathbf{x}, \lambda) F_\Gamma^+(\mathbf{x}, \lambda) \right]_{\lambda=0} = \tag{15}$$

$$\mathbb{E}_\mathbf{x}\left[ \mathbb{1}_A(\mathbf{X}(\tau_{A\cup B}^-)) \mathbb{1}_B(\mathbf{X}(\tau_{A\cup B}^+)) \left( \int_{\tau_{A\cup B}^-}^{\tau_{A\cup B}^+} \Gamma(\mathbf{X}(r))\, dr \right)^k \right]$$

where, again, the expectation is restricted to paths crossing through $\mathbf{x}$. Setting $k = 0$, this is simply $q_B^+(\mathbf{x}) q^-(\mathbf{x})$, the probability of an observed snapshot $\mathbf{x}$ being part of a transition path. With $k \geq 1$, it is natural to condition on snapshots being reactive by dividing by $q_A^-(\mathbf{x}) q_B^+(\mathbf{x})$.

$$\frac{\partial_\lambda^k \left[ F_\Gamma^+(\mathbf{x}; \lambda) F_\Gamma^-(\mathbf{x}; \lambda) \right]_{\lambda=0}}{q_A^-(\mathbf{x}) q_B^+(\mathbf{x})} \tag{16}$$

$$= \mathbb{E}_\mathbf{x}\left[ \left( \int_{\tau_{A\cup B}^-}^{\tau_{A\cup B}^+} \Gamma(\mathbf{X}(r))\, dr \right)^k \middle| \mathbf{X}(\tau_{A\cup B}^-) \in A, \mathbf{X}(\tau_{A\cup B}^+) \in B \right]$$

$$= \mathbb{E}_\mathbf{x}\left[ \left( \int_{\tau_{A\cup B}^-}^{\tau_{A\cup B}^+} \Gamma(\mathbf{X}(r))\, dr \right)^k \middle| A \to B \right] \text{ (abbreviation)}$$

Everything we say about transition paths stems originally from the functions $F_\Gamma^+$ and $F_\Gamma^-$ for various $\Gamma$, as well as the steady-state distribution $\pi$. Thus, we will now express the quantities of interest in terms of $\pi$, $F_\Gamma^+$, $F_\Gamma^-$, and their $\lambda$-derivatives. Section 2 will then explain how to compute them using short simulation data.

### 1.1.2 Reactive snapshot averages

In this subsection, we use committors to define statistics over the transition path ensemble at the level of snapshots, or equivalently, averages with respect to the reactive densities $\pi_{AB}$ and $\pi_{BA}$. The following subsection does the same at the level of paths, and uses the more general $F_\Gamma^\pm$ forecasts rather than just committors. The key to transforming forecasts into ensemble averages (at either level) is the *ergodic assumption*, which goes as follows. Let $\mathbf{Y}(t)$ denote all the hidden variables of the system responsible for apparent randomness, such as unresolved turbulence, so that the joint process $(\mathbf{X}(t), \mathbf{Y}(t))$ is completely deterministic. The hypothesis then states that there exists a probability density $p(\mathbf{x}, \mathbf{y})$ such that for any bounded observable $G(\mathbf{X}(t), \mathbf{Y}(t))$,

$$\lim_{T \to \infty} \frac{1}{T} \int_0^T G(\mathbf{X}(t), \mathbf{Y}(t))\, dt \tag{17}$$

$$= \int \left( \int G(\mathbf{x}, \mathbf{y}) p(\mathbf{x}, \mathbf{y})\, d\mathbf{y} \right) \pi(\mathbf{x})\, d\mathbf{x}$$

$$=: \int \Gamma(\mathbf{x}) \pi(\mathbf{x})\, d\mathbf{x} =: \langle \Gamma \rangle_\pi. \tag{18}$$

As a prototypical example, define

$$G(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & (\mathbf{x}, \mathbf{y}) \text{ comes from } A \text{ and goes to } B \\ 0 & \text{otherwise} \end{cases} \tag{19}$$

where knowledge of $\mathbf{y}$ lets us run the system forward and backward deterministically to evaluate the source and destination of $(\mathbf{x}, \mathbf{y})$. In this case, Eq. (17) becomes

$$\lim_{T \to \infty} \frac{1}{T} \int_0^T \mathbb{1}_A(\mathbf{X}(\tau_{A\cup B}^-(t))) \mathbb{1}_B(\mathbf{X}(\tau_{A\cup B}^+(t))\, dt$$

$$= \int q_A^-(\mathbf{x}) q_B^+(\mathbf{x}) \pi(\mathbf{x})\, d\mathbf{x} =: \langle q_A^- q_B^+ \rangle_\pi. \tag{20}$$

3

The left-hand side is the time fraction spent on the way from $A$ to $B$, and can be estimated from DNS. On the right hand side, $\Gamma(\mathbf{x}) = q_A^+(\mathbf{x})q_B^+(\mathbf{x})$, which we have already solved for using DGA. Substituting different combinations of $A$ and $B$ in Eq. (20) gives us the time fraction spend in the phases $B \to A$ ($\Gamma = q_B^- q_A^+$), $A \to A$ ($\Gamma = q_A^- q_A^+$), and $B \to B$ ($\Gamma = q_B^- q_B^+$). Fig. 1a shows these estimates from DNS and DGA, which is further described in section 2 below.

### 1.1.3 Transition path averages and currents

Thinking about transition paths from start to finish as discrete, coherent objects unlocks a richer description of the rare event, which is ultimately related to the reactive currents. In this section we use forecast functions (12) and (13) to express the transition rate and statistics of the path integrals (7) with a single framework centered around the *generalized rate*:

$$R_\Gamma(\lambda) := \lim_{T \to \infty} \frac{1}{T} \sum_{m=1}^{M_T} \exp\left(\lambda \int_{\tau_m^-}^{\tau_m^+} \Gamma(\mathbf{X}(r))\, dr\right) \tag{21}$$

The notation emphasizes that $R_\Gamma$ depends on the observable $\Gamma$ and the real parameter $\lambda$. To unpack this formula, first set $\lambda = 0$ and observe that $R_\Gamma(0) = \frac{M_T}{T}$ is the number of transitions per unit time, or rate, whose inverse is the average period of the full SSW life cycle. This is not to be confused with the asymmetric forward and backward rates,

$$k_{AB} = \frac{R_\Gamma(0)}{\langle q_A^- \rangle_\pi}, \qquad\qquad k_{BA} = \frac{R_\Gamma(0)}{\langle q_B^- \rangle_\pi} \tag{22}$$

which distinguish the $A \to B$ and $B \to A$ directions by how fast they occur. The factor $\langle q_A^- \rangle_\pi$ is the time fraction spent *having last been in $A$* rather than $B$, and $\langle q_B^- \rangle$ is the opposite. For example, if $A$ were very stable and $B$ very unstable, the system would spend most of its time in the basin of attraction of $A$, making $\langle q_A^- \rangle_\pi$ large and $k_{AB} \ll k_{BA}$. Asymmetric rates (or "rate constants") are very important for chemistry applications, but the symmetric rate is more useful to us presently.

All of these rates have been studied extensively with TPT and preceding theories. A novel idea that we introduce here is to include the exponential factor $\exp\left(\lambda \int \Gamma(\mathbf{X}(t))\, dt\right)$ to additionally study transition path integrals, though we do not present these results in this paper. The theoretical development below therefore reduces to classical TPT by replacing $\Gamma$ with 0.

Returning to (21), we divide through by $R_\Gamma(0)$:

$$\frac{R_\Gamma(\lambda)}{R_\Gamma(0)} = \lim_{T \to \infty} \frac{1}{M_T} \sum_{m=1}^{M_T} \exp\left(\lambda \int_{\tau_m^-}^{\tau_m^+} \Gamma(\mathbf{X}(r))\, dr\right)$$

$$= \mathbb{E}_{\text{paths}}\left[\exp\left(\lambda \int_{\tau_A^-}^{\tau_B^+} \Gamma(\mathbf{X}(r))\, dr\right)\right] \tag{23}$$

where the subscript "paths" distinguishes the expectation as over *all* transition paths, not just those crossing through a fixed $\mathbf{x}$ as in (16). The right side of (23) is a moment-generating function for the transition path integral (7). Differentiating in $\lambda$ yields the moments of that distribution, including its variance, skew, and kurtosis:

$$\frac{\partial_\lambda^k R_\Gamma(0)}{R_\Gamma(0)} = \mathbb{E}_{\text{paths}}\left[\left(\int_{\tau_A^-}^{\tau_B^+} \Gamma(\mathbf{X}(r))\, dr\right)^k\right], \tag{24}$$

Thus, $R_\Gamma(\lambda)$ contains much information about the transition ensemble as measured by path integrals.

We now express $R_\Gamma$ in terms of the forecast functions $F_\Gamma^+$ and $F_\Gamma^-$, again using the key assumption of ergodicity. We must convert Equation (21), a sum over transition paths $\sum_{m=1}^{M_T}(\cdot)$, into an integral over time $\int_{-T}^{T}(\cdot)\, dt$ and then (by ergodicity) into an integral over space $\int_{\mathbb{R}^d}(\cdot)\pi(\mathbf{x})\, d\mathbf{x}$. This approach extends the rate derivation in [612006Vanden-Eijnden] and [592021Strahan et al.Strahan, Antoszewski, Lorpaiboon, Vani, Weare,, and Dinner] to generalized rates.

To write the rate as a time integral, we introduce a dividing surface between $A$ and $B$ (such as a committor level surface) and use the fact that a transition path crosses such a surface an odd number of times. A mask is applied to the time integral to select only the time segments when a reactive trajectory segment is crossing this surface (+1 for positive crossings and −1 for negative crossings), resulting in unit weight for each transition path. To be more explicit, let $S$ be a region of state space that contains $A$ and excludes $B$, so that its boundary $C = \partial S$ is a dividing surface between $A$ and $B$. The generalized rate (21) can

then be written as the following time integral:

$$R_\Gamma(\lambda) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \lim_{T \to \infty} \frac{1}{T} \int_0^T \tag{25}$$

$$\exp\left(\lambda \int_{\tau^-_{A \cup B}(t)}^{\tau^+_{A \cup B}(t+\Delta t)} \Gamma(\mathbf{X}(r))\,dr\right) \times \tag{26}$$

$$\mathbb{1}_A(\mathbf{X}(\tau^-_{A \cup B}(t)))\,\mathbb{1}_B(\mathbf{X}(\tau^+_{A \cup B}(t+\Delta t))) \times \tag{27}$$

$$\left[\mathbb{1}_S(\mathbf{X}(t))\,\mathbb{1}_{S^c}(\mathbf{X}(t+\Delta t))\right. \tag{28}$$

$$\left. -\,\mathbb{1}_{S^c}(\mathbf{X}(t))\,\mathbb{1}_S(\mathbf{X}(t+\Delta t))\right]dt \tag{29}$$

The idea is to restrict the interval $(0,T)$ to the collection of time intervals $(t, t+\Delta t)$ during which the path crosses the surface $\partial S$. Line (27) applies a mask picking out transition path segments, which are those that come from $A$ and next go to $B$. Line (28) applies a further mask picking out the narrow time intervals when $\mathbf{X}(t)$ exits the region from $S$ to $S^c$, while line (29) subtracts the backward crossings from $S^c$ to $S$. Using ergodicity, we can replace the time integral with a space integral and insert conditional expectations inside. For example, the part of the integrand

$$\exp\left(\lambda \int_{t+\Delta t}^{\tau^+_{A \cup B}(t+\Delta t)} \Gamma(\mathbf{X}(r))\,dr\right) \times \tag{30}$$
$$\mathbb{1}_B(\mathbf{X}(\tau^+_{A \cup B}(t+\Delta t)))\,\mathbb{1}_{S^c}(\mathbf{X}(t+\Delta t))$$

becomes, after taking conditional expectations,

$$\mathbb{E}\left[\mathbb{1}_{S^c}(\mathbf{X}(t+\Delta t))F_\Gamma^+(\mathbf{X}(t+\Delta t))\big|\mathbf{X}(t) = \mathbf{x}\right] \tag{31}$$
$$=: \mathcal{T}^{\Delta t}\left[\mathbb{1}_{S^c}F_\Gamma^+\right](\mathbf{x})$$

Where the *transition operator* is defined as $\mathcal{T}^{\Delta t} f(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[f(\mathbf{X}(\Delta t))]$. Applying similar logic to all terms in the integrand, we have the following generalized rate formula:

$$R_\Gamma(\lambda) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \int_{\mathbb{R}^d} F_\Gamma^-(\mathbf{x}; \lambda) \times \tag{32}$$
$$\left\{\mathbb{1}_S \mathcal{T}^{\Delta t}\left[\mathbb{1}_{S^c}F_\Gamma^+\right] - \mathbb{1}_{S^c}\mathcal{T}^{\Delta t}\left[\mathbb{1}_S F_\Gamma^+\right]\right\}(\mathbf{x})\pi(\mathbf{x})\,d\mathbf{x}$$

which holds for any $S$ enclosing $A$ and disjoint from $B$. This is a form estimable from short simulation data, which the next section will explain.

The rate formula (32) is suggestive of a surface integral, counting hopping events across the surface $\partial S$. In fact, the reactive current $\mathbf{J}_{AB}$ is defined as the vector field whose surface integral is equal to the symmetric rate:

$$R_\Gamma(0) = \int_C \mathbf{J}_{AB} \cdot \mathbf{n}\,d\sigma \tag{33}$$

We have visualized $\mathbf{J}_{AB}$ in section 3 using a discretization of the integrand in (32). Note that the integral does *not* depend on the specific dividing surface we choose, which implies that $\mathbf{J}_{AB}$ is divergence-free outside of $A \cup B$, but has a source of field lines at $A$ and a sink at $B$. Since every dividing surface supports the same total flux, large local current magnitude means a constrained reaction mechanism, as demonstrated by the early transition states (larger values of $U(30\text{ km})$) in Fig. 6a.

We have now completely described the mathematics of TPT, and our extensions to it. Exact knowledge of $\pi$, $F_\Gamma^+$, $F_\Gamma^-$, and their $\lambda$-derivatives is enough to generate all of the figures shown so far. The next section explains both how to compute these fundamental ingredients from data and assemble them into generalized rates.

# 2 Numerical method: dynamical Galerkin approximation (DGA)

## 2.1 Feynman-Kac formulae

We now sketch the numerical method, following [602019Thiede et al.Thiede, Giannakis, Dinner,, and Weare, 592021Strahan et al.Strahan, An and [202021Finkel et al.Finkel, Webber, Abbot, Gerber,, and Weare]. Equations (12) and (13) involve an integral in time all the way from $t = 0$ to $t = \tau^+_{A \cup B}$ or (in backward time) to $\tau^-_{A \cup B}$, when $\mathbf{X}(t)$ hits either $A$ or $B$ after wandering through state space

for an indeterminate period. This would seem to require long trajectories to estimate. Section 4 of the main paper explained heuristically how short trajectories can be substituted. Here, we present DGA more generally, writing $F_\Gamma^\pm$ as solutions to partial differential equations called Feynman-Kac formulae [472003Oksendal, **?**, **?**], which read

$$\begin{cases} (\mathcal{L} + \lambda\Gamma(\mathbf{x}))F_\Gamma^+(\mathbf{x};\lambda) = 0 & \mathbf{x} \in D \\ F_\Gamma^+(\mathbf{x}) = 0 & \mathbf{x} \in A \\ F_\Gamma^+(\mathbf{x}) = 1 & \mathbf{x} \in B \end{cases} \tag{34}$$

$$\text{where } \mathcal{L}\phi(\mathbf{x}) := \lim_{\Delta t \to 0} \frac{\mathbb{E}_\mathbf{x}[\phi(\mathbf{X}(\Delta t))] - \phi(\mathbf{x})}{\Delta t} \tag{35}$$

$$\begin{cases} (\widetilde{\mathcal{L}} + \lambda\Gamma(\mathbf{x}))F_\Gamma^-(\mathbf{x};\lambda) = 0 & \mathbf{x} \in D \\ F_\Gamma^-(\mathbf{x}) = 1 & \mathbf{x} \in A \\ F_\Gamma^-(\mathbf{x}) = 0 & \mathbf{x} \in B \end{cases} \tag{36}$$

$$\text{where } \widetilde{\mathcal{L}}\phi(\mathbf{x}) := \lim_{\Delta t \to 0} \frac{\mathbb{E}_\mathbf{x}[\phi(\mathbf{X}(-\Delta t))] - \phi(\mathbf{x})}{\Delta t} \tag{37}$$

The linear operators $\mathcal{L}$ and $\widetilde{\mathcal{L}}$ are known as the forward and backward infinitesimal generators, pushing observable functions $\phi$ forward or backward in time analogously to a material derivative in fluid mechanics. The first term in the numerator of (35) is the transition operator. The backward-in-time expectations are defined specifically for the *equilibrium* process, leading to $\widetilde{\mathcal{L}}\phi(\mathbf{x}) = \frac{1}{\pi(\mathbf{x})}\mathcal{L}^*[\pi\phi](\mathbf{x})$, where $\mathcal{L}^*$ is the adjoint of $\mathcal{L}$ with respect to the reference (Lebesgue) measure $d\mathbf{x}$ (see discussion following Eq. (28) of the main text). Equivalently, $\widetilde{\mathcal{L}}$ is the adjoint of $\mathcal{L}$ with respect to the steady-state measure $\pi(\mathbf{x})\,d\mathbf{x}$. In addition, we have the stationary Fokker-Planck equation for $\pi$ itself:

$$\begin{cases} \mathcal{L}^*\pi(\mathbf{x}) = 0 & \mathbf{x} \in \mathbb{R}^d \\ \int_{\mathbb{R}^d} \pi(\mathbf{x})\,d\mathbf{x} = 1 \end{cases} \tag{38}$$

We can further obtain equations for the derivatives of $F_\Gamma^\pm$ with respect to $\lambda$, using the Kac moment method [**?**] as also described in [202021Finkel et al.Finkel, Webber, Abbot, Gerber,, and Weare]. Differentiating the equation (35) itself in $\lambda$ and setting $\lambda = 0$ yields a recursive sequence of equations for the higher derivatives of $F_\Gamma^+$:

$$\mathcal{L}[\partial_\lambda^k F_\Gamma^+](\mathbf{x};0) = -k\Gamma\partial_\lambda^{k-1}F_\Gamma^+(\mathbf{x};0), \qquad\qquad k \geq 1, \tag{39}$$

with boundary conditions $\partial_\lambda^k F_\Gamma^+\big|_{A \cup B} = 0$. The same procedure can be applied to the aftcast $F_\Gamma^-$. Thus, our entire numerical pipeline boils down to solving equations of the form (34), (36), and (38), as well as the inhomogeneous version (39).

## 2.2 Discretization of Feynman-Kac formulae

We will now describe how to discretize and solve these three equations, which requires three similar but distinct procedures.

First we attack (34). The generator of a diffusion processes can be expressed as a partial differential operator, and so the above equations are PDEs over state space. PDEs cannot be practically discretized in high dimensions, but the essential property of spatial locality allows for data-driven approximation with short trajectories, using the probabilistic definition in (35) and (37). This is how we use our large data set of short trajectories,

$$\{\mathbf{X}_n(t) : 0 \leq t \leq \Delta t\}_{n=1}^N, \tag{40}$$

where the initial points $\mathbf{X}_n(0)$ are drawn from a *sampling measure* $\mu$, which we will define in the following subsection. To discretize Eq. (34), we first eliminate the numerically problematic limit $\Delta t \to 0$ and integrate the equation using Dynkin's Formula [472003Oksendal, **?**]: for any stopping time $\theta > 0$,

$$\mathbb{E}_\mathbf{x}[f(\mathbf{X}(\theta))] = f(\mathbf{x}) + \mathbb{E}_\mathbf{x}\left[\int_0^\theta \mathcal{L}f(\mathbf{X}(t))\,dt\right] \tag{41}$$

Here we take $\theta = \min(\Delta t, \tau_{A \cup B})$. In other words, we artificially halt the $n$th trajectory $\mathbf{X}_n(t)$ if it wanders into $A \cup B$ before the terminal time $\Delta t$. The $n$th stopping time from the data set is called $\theta_n$. The operator on the left-hand side of (41) is known as the *stopped* transition operator $\mathcal{T}^\theta$. Applying it to the unknown forecast function $F_\Gamma^+$ in (34) and using the fact $\mathcal{L}F_\Gamma^+ = -\lambda\Gamma F_\Gamma^+$, we get

$$\mathcal{T}^\theta F_\Gamma^+(\mathbf{x};\lambda) = F_\Gamma^+(\mathbf{x};\lambda) - \lambda\mathbb{E}_\mathbf{x}\left[\int_0^\theta \Gamma(\mathbf{X}(t))F_\Gamma^+(\mathbf{X}(t);\lambda)\,dt\right] \tag{42}$$

6

To be more concise, we define the integral operator $\mathcal{K}_\Gamma^\theta f(\mathbf{x}) = \mathbb{E}_\mathbf{x}\left[\int_0^\theta \Gamma(\mathbf{X}(t))f(\mathbf{X}(t))\,dt\right]$, and write an integrated version of (34):

$$
\begin{cases}
(\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)F_\Gamma^+(\mathbf{x};\lambda) = 0 & \mathbf{x} \in D \\
F_\Gamma^+(\mathbf{x};\lambda) = 0 & \mathbf{x} \in A \\
F_\Gamma^+(\mathbf{x};\lambda) = 1 & \mathbf{x} \in B
\end{cases}
\tag{43}
$$

To discretize this equation and impose regularity on the solution, we approximate $F_\Gamma^+$ as a finite linear combination with coefficients $w_j(F_\Gamma^+(\mathbf{x};\lambda))$, which we abbreviate $w_j(\lambda)$ for simplicity:

$$
F_\Gamma^+(\mathbf{x};\lambda) \approx \hat{F}_\Gamma^+(\mathbf{x};\lambda) + \sum_{j=1}^M w_j(\lambda)\phi_j(\mathbf{x};\lambda)
\tag{44}
$$

where $\hat{F}_\Gamma^+$ is a guess function obeying the boundary conditions on $A \cup B$, and $\{\phi_j\}_{j=1}^M$ is a collection of basis functions that are zero on $A \cup B$, which will be defined in the following subsection. The task is now to solve for the coefficients $w_j(\lambda)$. Equation (43) becomes a system of linear equations in $w_j(\lambda)$:

$$
\sum_{j=1}^M w_j(\lambda)(\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)\phi_j(\mathbf{x};\lambda) = -(\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)\hat{F}_\Gamma^+(\mathbf{x};\lambda)
\tag{45}
$$

Since the transfer and integral operators are expectations over the future state of the system beginning at $\mathbf{x}$, we can estimate their action at $\mathbf{x} = \mathbf{X}_n(0)$ (a short-trajectory starting point) as

$$
(\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)\phi_j(\mathbf{X}_n(0);\lambda) \approx \phi_j(\mathbf{X}_n(\theta_n);\lambda) - \phi_j(\mathbf{X}_n(0);\lambda)
\tag{46}
$$
$$
+ \lambda \int_0^{\theta_n} \Gamma(\mathbf{X}_n(t))\phi_j(\mathbf{X}_n(t);\lambda)\,dt
$$

or, if multiple independent trajectories are launched from $\mathbf{x}$, we can average over them. Applying this to every short trajectory and plugging into Eq. (45), we obtain a system of $N$ equations in $M$ unknowns. In practice, $N \gg M$, meaning we have many more trajectories than basis functions, and the system is overdetermined. A unique, and regularized, solution is obtained by casting it into weak form: we multiply both sides by $\phi_i(\mathbf{x})$ and integrate over state space:

$$
\sum_{j=1}^M w_j(\lambda)\left\langle\phi_i, (\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)\phi_j\right\rangle_\zeta = -\left\langle\phi_i, (\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)\hat{F}_\Gamma^+\right\rangle_\zeta
\tag{47}
$$

where the inner products are defined with respect to a measure $\zeta$:

$$
\langle f, g\rangle_\zeta = \int f(\mathbf{x})g(\mathbf{x})\zeta(\mathbf{x})\,d\mathbf{x}
\tag{48}
$$

With our finite data set, we approximate the inner product by a sum over pairs of points. Given that $\mathbf{X}_n(0) \sim \mu$, the law of large numbers ensures that for any bounded function $H(\mathbf{x})$,

$$
\frac{1}{N}\sum_{n=1}^N H(\mathbf{X}_n(0)) \approx \int H(\mathbf{x})\mu(\mathbf{x})\,d\mathbf{x}
\tag{49}
$$

becomes more accurate as $N \to \infty$. Thus we set $H(\mathbf{x}) = \phi_i(\mathbf{x})(\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)\phi_j(\mathbf{x})$ as estimated by (46), approximate the inner products with $\zeta = \mu$, plug them into (47), and solve the $M \times M$ system of linear equations for $w_j(\lambda)$.

Next we address (38). The integrated version of (38) is found by observing that for any bounded function $H$,

$$
\mathbb{E}_{\mathbf{X}(0)\sim\pi}[H(\mathbf{X}(\Delta t))] = \mathbb{E}_{\mathbf{X}(\Delta t)\sim\pi}[H(\mathbf{X}(\Delta t))]
\tag{50}
$$

where $\mathbf{X}(0) \sim \pi$ means the initial condition is drawn from equilibrium, and thus so is $\mathbf{X}(\Delta t)$ since $\pi$ is stationary. Of course, our initial data is *not* distributed according to $\pi$, but rather by $\mu$; the goal is to solve for the *change of measure* $\frac{d\pi}{d\mu}(\mathbf{x})$. Writing (50)

as an integral,

$$\int \mathcal{T}^{\Delta t} H(\mathbf{x}) \pi(\mathbf{x}) \, d\mathbf{x} = \int H(\mathbf{x}) \pi(\mathbf{x}) \, d\mathbf{x} \tag{51}$$

$$0 = \int (\mathcal{T}^{\Delta t} - 1) H(\mathbf{x}) \pi(\mathbf{x}) \, d\mathbf{x} \tag{52}$$

$$= \int (\mathcal{T}^{\Delta t} - 1) H(\mathbf{x}) \frac{d\pi}{d\mu}(\mathbf{x}) \mu(\mathbf{x}) \, d\mathbf{x} \tag{53}$$

$$= \left\langle (\mathcal{T}^{\Delta t} - 1) H, \frac{d\pi}{d\mu} \right\rangle_{\mu} \tag{54}$$

As this holds for every bounded $H$, we enforce the equation for $H = \phi_i$, $i = 1, \ldots, M$ and approximate $\frac{d\pi}{d\mu} = \sum_{j=1}^{M} w_j \left( \frac{d\pi}{d\mu} \right) \phi_j$, resulting in a homogeneous linear system for the coefficients $w_j$ similar to (47). The matrix elements are

$$\langle (\mathcal{T}^{\Delta t} - 1) \phi_i, \phi_j \rangle_{\mu} = \langle \phi_i, (\mathcal{T}^{\Delta t} - 1)^*_{\mu} \phi_j \rangle_{\mu} \tag{55}$$

where $(\cdot)^*_{\mu}$ denotes the adjoint operator with respect to $\mu$. For the indicator basis that we use, these inner products yield the entries of the Markov matrix $P_{ij}$ described in the main text. [If we were to divide by $\Delta t$ and take the limit $\Delta t \to 0$, we would recover the strong form of the Fokker-Planck equation, (38).] We solve this homogeneous system by $QR$ decomposition. Note that there are no boundary conditions, and the trajectories need not be stopped early. Instead there is a normalization condition, which we enforce as $\sum_{n=1}^{N} \frac{d\pi}{d\mu}(\mathbf{X}_n(0)) = 1$. Furthermore, the basis $\{\phi_j\}$ is different between $\frac{d\pi}{d\mu}$ and $F_{\Gamma}^+$. To ensure that the matrix has a nontrivial null vector, one can add a constant vector to the basis. However, the basis of indicators that we use guarantees a null space automatically. Given the weights $\frac{d\pi}{d\mu}(\mathbf{X}_n(0))$, we can take any ergodic average $\Gamma$ by inserting the change of measure:

$$\langle \Gamma \rangle_{\pi} = \int_{\mathbb{R}^d} \Gamma(\mathbf{x}) \pi(\mathbf{x}) \, d\mathbf{x} = \int_{\mathbb{R}^d} \Gamma(\mathbf{x}) \, d\mathbf{x} \approx \sum_{n=1}^{N} \Gamma(\mathbf{X}_n(0)) \frac{d\pi}{d\mu}(\mathbf{X}_n(0)) \tag{56}$$

For the specific case of a Markov state model, we can estimate the $\pi$-weighted state-space integral by decomposing it over clusters:

$$\langle \Gamma \rangle_{\pi} = \sum_{j=1}^{M} \int_{S_j} \Gamma(\mathbf{x}) \pi(\mathbf{x}) \, d\mathbf{x} \tag{57}$$

$$\approx \sum_{j=1}^{M} w_j(\pi) \times (\text{average of } \Gamma \text{ over } S_j) \tag{58}$$

$$\approx \sum_{j=1}^{M} w_j(\pi) \frac{\sum_{n=1}^{N} \Gamma(\mathbf{X}_n(0)) \mathbb{1}_{S_j}(\mathbf{X}_n(0))}{\sum_{n'=1}^{N} \mathbb{1}_{S_j}(\mathbf{X}_{n'}(0))} \tag{59}$$

$$= \sum_{n=1}^{N} \Gamma(\mathbf{X}_n(0)) \frac{w_{j(n)}(\pi)}{\#\{n' : j(n') = j(n)\}} \tag{60}$$

where $j(n)$ is defined as the cluster that $\mathbf{X}_n(0)$ is assigned to, i.e., $\mathbf{X}_n(0) \in S_{j(n)}$. Hence the change of measure for a Markov state model is the coefficient of $\Gamma(\mathbf{X}_n(0))$ in the last equation. Note that it sums to one over all data points, as it must.

Finally, we address the time-reversed Kolmogorov equation (36). The only modification from (34) is that the inner products in (47) are interpreted in backward time, i.e., with all trajectories reversed, $\mathbf{X}_n(\Delta t)$ becoming the beginning and $\mathbf{X}_n(0)$ becoming the end. The problem is that $\mathbf{X}_n(\Delta t)$ is not distributed according to $\mu$, and so we cannot use the same Monte Carlo inner product as in Eq. (49) with reference measure $\zeta = \pi$. However, we can solve the problem by reweighting with the change of measure as follows, leading to $\zeta = \pi$. We let the trajectory be discrete in time, i.e.,

$$\mathbf{X}_n = \left[ \mathbf{X}_n(0), \mathbf{X}_n\left(\frac{\Delta t}{K}\right), \mathbf{X}_n\left(\frac{2\Delta t}{K}\right), \ldots, \mathbf{X}_n(\Delta t) \right] \tag{61}$$

and consider functionals $H[\mathbf{X}_n]$ of the whole trajectory. Defining the transition density $p(\mathbf{x}, \mathbf{y})$ for each step of size $\Delta t$, the expectation of $H$ with $\mathbf{X}_n(0) \sim \pi$ is given by

$$\mathbb{E}_{\mathbf{X}(0) \sim \pi} H[\mathbf{X}] = \int d\mathbf{x}_0 \pi(\mathbf{x}_0) \int d\mathbf{x}_1 p(\mathbf{x}_0, \mathbf{x}_1) \int \ldots \int d\mathbf{x}_K p(\mathbf{x}_{K-1}, \mathbf{x}_K) H[\mathbf{x}_0, \ldots, \mathbf{x}_K] \tag{62}$$

The time reversal step explicitly assumes the *equilibrium* backward process, leading to a backward transition kernel $\widetilde{p}(\mathbf{y},\mathbf{x}) = \frac{\pi(\mathbf{x})}{\pi(\mathbf{y})}p(\mathbf{x},\mathbf{y})$. Inserting this throughout converts the expectation over $\mathbf{X}(0)$ into an expectation over $\mathbf{X}(\Delta t)$:

$$\mathbb{E}_{\mathbf{X}(0)\sim\pi}H[\mathbf{X}] = \int d\mathbf{x}_0 \pi(\mathbf{x}_0) \int d\mathbf{x}_1 \frac{\pi(\mathbf{x}_1)}{\pi(\mathbf{x}_0)}\widetilde{p}(\mathbf{x}_1,\mathbf{x}_0) \int \tag{63}$$

$$\ldots \int d\mathbf{x}_K \frac{\pi(\mathbf{x}_K)}{\pi(\mathbf{x}_{K-1})}\widetilde{p}(\mathbf{x}_K,\mathbf{x}_{K-1})H[\mathbf{x}_0,\ldots,\mathbf{x}_K] \tag{64}$$

$$= \int d\mathbf{x}_K \pi(\mathbf{x}_K) \int d\mathbf{x}_{K-1}\widetilde{p}(\mathbf{x}_K,\mathbf{x}_{K-1}) \int \ldots \int d\mathbf{x}_0 \widetilde{p}(\mathbf{x}_1,\mathbf{x}_0)H[\mathbf{x}_0,\ldots,\mathbf{x}_K] \tag{65}$$

$$= \widetilde{\mathbb{E}}_{\mathbf{X}(\Delta t)\sim\pi}H[\mathbf{X}] \tag{66}$$

where $\widetilde{\mathbb{E}}$ denotes backward-in-time expectation. This is precisely what we need to apply (49) to the time-reversed process, namely, define $H$ such that

$$\phi_i(\mathbf{X}(\Delta t))(\mathcal{T}^\theta - 1 + \lambda\mathcal{K}_\Gamma^\theta)\phi_j(\mathbf{X}(\Delta t)) =: H[\mathbf{X}] \tag{67}$$

and then integrate over state space weighted by $\pi$, turning the left-hand side into an inner product:

$$\langle \phi_i, (\widetilde{\mathcal{T}}^\theta - 1 + \lambda\widetilde{\mathcal{K}}_\Gamma^\theta)\phi_j\rangle_\pi = \widetilde{\mathbb{E}}_{\mathbf{X}(\Delta t)\sim\pi}H[\mathbf{X}] \tag{68}$$

$$= \mathbb{E}_{\mathbf{X}(0)\sim\pi}H[\mathbf{X}] \approx \sum_{n=1}^{N} H[\mathbf{X}_n]\frac{d\pi}{d\mu}(\mathbf{X}_n(0)) \tag{69}$$

The right-hand side of Eq. (49) can be estimated similarly, also with $\zeta = \pi$.

In both forward- and backward-time estimates, we never solve for $F_\Gamma^+(\mathbf{x};\lambda)$ or $F_\Gamma^-(\mathbf{x};\lambda)$ with nonzero $\lambda$; rather, we repeat the recursion process with Eq. (39). This is equivalent to implicitly differentiating the discretized system Eq. (47).

## 2.3 Rate estimate and numerical benchmarking

To estimate generalized rates (in particular, the ordinary rate), we reproduce here the rate estimate from [592021Strahan et al.Strahan, Antoszew for reference, which is an almost-direct implementation of the formula (32), repeated here:

$$R_\Gamma(\lambda) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \int_{\mathbb{R}^d} F_\Gamma^-(\mathbf{x};\lambda) \times \tag{70}$$

$$\left\{ \mathbb{1}_S \mathcal{T}^{\Delta t}\left[\mathbb{1}_{S^c}F_\Gamma^+\right] - \mathbb{1}_{S^c}\mathcal{T}^{\Delta t}\left[\mathbb{1}_S F_\Gamma^+\right]\right\}(\mathbf{x})\pi(\mathbf{x})\,d\mathbf{x}$$

In principle, the integral could be estimated directly with any choice of dividing surface $S$, but the sum would only use data either exiting $S$ (first term) or entering $S$ (second term). We can use all the data at once and improve numerical stability by averaging over multiple such surfaces, and furthermore converting the transition operator $\mathcal{T}^{\Delta t}$ into the generator $\mathcal{L}$. However, we cannot simply take the limit under the integral due to the discontinuity in $\mathbb{1}_S$. Instead we get a smooth function into the integrand with the following steps. First, replace $\mathbb{1}_S$ with $1 - \mathbb{1}_{S^c}$ everywhere:

$$R_\Gamma(\lambda) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \int_{\mathbb{R}^d} F_\Gamma^-(\mathbf{x};\lambda) \times \tag{71}$$

$$\left\{(1 - \mathbb{1}_{S^c})\mathcal{T}^{\Delta t}\left[\mathbb{1}_{S^c}F_\Gamma^+\right] - \mathbb{1}_{S^c}\mathcal{T}^{\Delta t}\left[(1-\mathbb{1}_{S^c})F_\Gamma^+\right]\right\}(\mathbf{x})\pi(\mathbf{x})\,d\mathbf{x} \tag{72}$$

$$= \lim_{\Delta t \to 0}\frac{1}{\Delta t}\int_{\mathbb{R}^d} F_\Gamma^-(\mathbf{x};\lambda)\left\{\mathcal{T}^{\Delta t}\left[\mathbb{1}_{S^c}F_\Gamma^+\right] - \mathbb{1}_{S^c}\mathcal{T}^{\Delta t}F_\Gamma^+\right\}(\mathbf{x})\pi(\mathbf{x})\,d\mathbf{x} \tag{73}$$

$$\tag{74}$$

Next, add and subtract $\mathbb{1}_{S^c}F_\Gamma^+$ inside the integrand.

$$R_\Gamma(\lambda) = \lim_{\Delta t \to 0}\int_{\mathbb{R}^d} F_\Gamma^-(\mathbf{x};\lambda)\left\{\frac{\mathcal{T}^{\Delta t} - 1}{\Delta t}\left[\mathbb{1}_{S^c}F_\Gamma^+\right] - \mathbb{1}_{S^c}\frac{\mathcal{T}^{\Delta t} - 1}{\Delta t}F_\Gamma^+\right\}(\mathbf{x})\pi(\mathbf{x})\,d\mathbf{x} \tag{75}$$

At this point it is tempting to take the limit inside the integral, as $(\mathcal{T}^{\Delta t} - 1)/\Delta t$ formally approaches $\mathcal{L}$. But the first term acts on a discontinuous function, which won't have a well-defined time derivative. We first replace $\mathbb{1}_{S^c}$ with a smooth function (on $D$), as follows.

Let $K : \mathbb{R}^d \to [0, 1]$ be a function that increases from 0 on set $A$ to 1 on set $B$ (for instance, the committor). Let $S_\zeta = \{\mathbf{x} : K(\mathbf{x}) \leq \zeta\}$ for $\zeta \in (0, 1)$, and integrate both sides over $\zeta$, noting that $\int_0^1 \mathbb{1}_{S_\zeta^c}(\mathbf{x}) \, d\zeta = \int_0^1 \mathbb{1}\{K(\mathbf{x}) > \zeta\} \, d\zeta = K(\mathbf{x})$.

$$\int_0^1 R_\Gamma(\lambda) \, d\zeta = \lim_{\Delta t \to 0} \int_{\mathbb{R}^d} F_\Gamma^-(\mathbf{x}; \lambda) \left\{ \frac{\mathcal{T}^{\Delta t} - 1}{\Delta t} [KF_\Gamma^+] - K\mathcal{L}F_\Gamma^+ \right\}(\mathbf{x})\pi(\mathbf{x}) \, d\mathbf{x} \tag{76}$$

Now we can move the limit inside and use the PDE to find

$$R_\Gamma(\lambda) = \int_{\mathbb{R}^d} F_\Gamma^-(\mathbf{x}; \lambda) \left\{ \mathcal{L}[KF_\Gamma^+](\mathbf{x}) + \lambda K(\mathbf{x})\Gamma(\mathbf{x})F_\Gamma^+(\mathbf{x}) \right\} \pi(\mathbf{x}) \, d\mathbf{x} \tag{77}$$

This formula can be estimated directly from knowledge of $F_\Gamma^-, F_\Gamma^+$, and $\pi$, using the ergodic assumption and with a discrete finite difference in time to estimate $\mathcal{L}[KF_\Gamma^+]$, i.e.,

$$\mathcal{L}[KF_\Gamma^+](\mathbf{X}_n(0)) \approx \frac{K(\mathbf{X}_n(\Delta t))F_\Gamma^+(\mathbf{X}_n(\Delta t)) - K(\mathbf{X}_n(0))F_\Gamma^+(\mathbf{X}_n(0))}{\Delta t} \tag{78}$$

Derivatives with respect to $\lambda$ can be found by iterating the product rule, as we have solved for the derivatives of $F_\Gamma^+$ and $F_\Gamma^-$.

To validate DGA numerically, we can compare to the results of DNS. In [202021Finkel et al.Finkel, Webber, Abbot, Gerber,, and Weare] (Fig. 7), we saw convergence of the DGA committor to the DNS committor across state space as sample size and lag time were increased. Here, we turn our attention to summary statistics of interest for full transition paths, not just forecasting. This will benchmark our current DGA implementation for comparison with future algorithmic developments.

Fig. 1a displays the time fractions spent in each phase of the SSW lifecycle: $A \to B$, $B \to A$, $A \to A$, and $B \to B$, including estimates from DNS (cyan) and DGA (red) and their uncertainties. The DGA estimate of the $A \to B$ time fraction is a $\pi$-weighted average of $q_A^-(\mathbf{x})q_B^+(\mathbf{x})$ over state space,

$$\langle q_A^- q_B^+ \rangle_\pi = \int q_A^-(\mathbf{x})q_B^+(\mathbf{x})\pi(\mathbf{x}) \, d\mathbf{x} \tag{79}$$

and similarly for the other phases (section 1 above justifies this formula rigorously, and section 2 above details the numerical computation of the integral). The DGA error bars are generated by repeating the entire pipeline three times with different short trajectory realizations. The bar height shows the mean, and the error bars show the minimum and maximum. The DNS error bars are generated by bootstrap resampling (with replacement) 500 times from the control simulation, treating an entire SSW lifecycle as a single unit (from the beginning of one $A \to B$ transition until the beginning of the next one). This assumes no memory between successive events, which we have found to be reasonable; there is insignificant autocorrelation between consecutive return periods. The bars extend two root-mean-squared errors in both directions, enclosing a 95% confidence interval. To first order, DGA agrees well with DNS on the fraction of time spent in each phase. $A$ is the more stable of the two regimes, accounting for $\sim 50\%$ of the time compared to the $\sim 40\%$ of time spent in the orbit of $B$. The transition events are both an order of magnitude shorter, with $B \to A$ taking slightly longer on average. DGA ranks the $A \to B$ and $B \to A$ time fractions correctly, despite a bias in the absolute magnitudes.

The numbers in Fig. 1a are only relative durations; they do not tell us how long a full life cycle takes. That number is given by (one over) the rate. Fig. 1b shows three different rate estimates (that is, the generalized rate with $\Gamma = 0$) using the formulas above. The cyan bars come from DNS, counting the number of $A \to B$ transitions per unit time. Of course, this equals the number of $B \to A$ transitions per unit time, so the $A \to B$ and $B \to A$ cyan bars are identical. Error bars come from bootstrapping, as with the relative durations. The red bars come from DGA, and these estimates are not technically symmetric. The DGA estimate labeled $A \to B$ integrates $\mathbf{J}_{AB} \cdot \mathbf{n}$ over dividing surfaces with $\mathbf{n}$ pointing away from $A$ toward $B$, while the estimate labeled $B \to A$ integrates $\mathbf{J}_{BA} \cdot \mathbf{n}$ over surfaces with $\mathbf{n}$ pointing away from $B$ toward $A$. Numerical and sampling errors cause slight differences between them, but Fig. 1b shows them both to come within 20% of the DNS estimate.

DGA estimates should converge with increasing $M$ (cluster number) and $N$ (short-trajectory ensemble size). Larger $M$ makes the approximation space $\{\phi_1, \ldots, \phi_M\}$ more expressive, making finer estimates possible. However, as $M$ grows, we need more short trajectories $N$ to robustly estimate the entries of the expanding matrix (26). Conversely, as $M$ shrinks, $P_{ij}$ will become closer to diagonal, because trajectories will escape from their starting cluster less frequently. Thus $\Delta t$ would have to increase when $M$ decreases. The optimal choice for a given model will depend on the relative costs of integrating the model, building basis sets, and solving large linear systems on different computer architectures. With our choice of $M = 1500$, increasing $N$ from $5 \times 10^4$ to $3 \times 10^6$ does not change the DGA point estimates very much, but shrinks the error bars by a factor of $\sim 4$. To further reduce the bias in Fig. 1, we would likely need more refined basis functions, perhaps using nonlinear features as input to K-means. For generalized rates such as transit time $\tau_B^+ - \tau_A^-$ and total heat flux $\int_{\tau_A^-}^{\tau_B^+} \overline{v'T'}(30 \text{ km}) \, dt$, a second-order calculation is required using Eq. (39), which causes errors to propagate further. The errors are not yet well-controlled enough to present the results of generalized rates. We do not yet have theoretical guarantees or optimal prescriptions for DGA parameters, but given the flexibility and parallelizability of the method, we believe it has much room for growth.
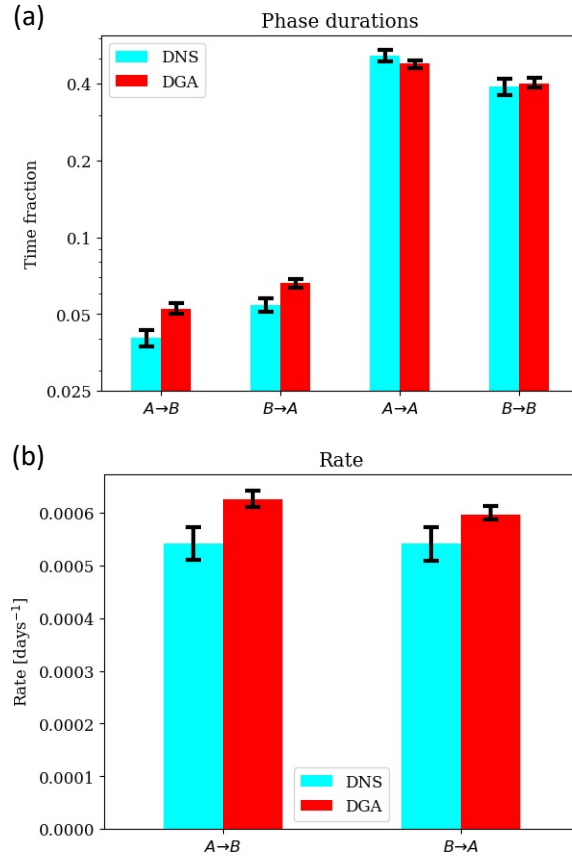
Figure 1: **DGA benchmarks and comparison to DNS.** (a) Time fractions spent in each phase. (b) Total SSW rate estimated using both $\mathbf{J}_{AB}$ and $\mathbf{J}_{BA}$; the two cyan columns, DNS estimates, are identical.

## 2.4 Visualization method

The two-dimensional projections in the main paper are generated with the following procedure. Let $\mathbf{y} = \mathbf{Y}(\mathbf{x})$ be an observable subspace, typically with dimension much less than that of $\mathbf{x}$ (usually two). Any scalar field $F(\mathbf{x})$, such as the committor, has a projection $F^{\mathbf{Y}}(\mathbf{y})$ onto this subspace by

$$F^{\mathbf{Y}}(\mathbf{y}) = \int F(\mathbf{x})\pi(\mathbf{x})\delta(\mathbf{Y}(\mathbf{x}) - \mathbf{y})\,d\mathbf{x} \tag{80}$$

In practice, the $\mathbf{y}$ space is partitioned into grid boxes $d\mathbf{y}$, and the integral is estimated from the dataset, yielding

$$F^{\mathbf{Y}}(\mathbf{y}) = \frac{1}{N}\sum_{n=1}^{N} F(\mathbf{X}_n(0))\frac{d\pi}{d\mu}(\mathbf{X}_n(0))\mathbb{1}_{d\mathbf{y}}(\mathbf{Y}(\mathbf{X}_n(0))) \tag{81}$$

where $\mathbb{1}_{d\mathbf{y}}(\mathbf{Y}(\mathbf{x})) = 1$ if $\mathbf{Y}(\mathbf{x}) \in d\mathbf{y}$ and zero otherwise. In words, we simply take a weighted average over all data points $\mathbf{X}_n(0)$ that project onto the grid box $d\mathbf{y}$, with weights given by the change of measure. In Fig. 4, we use $q_B^+$ and $q_A^-$ for $F$. In Fig. 5, we use $F = \pi$ (a), $F = q_A^- q_A^+$ (b), $F = q_A^- q_B^+$ (c), $F = q_B^- q_A^+$ (d), and $F = q_B^- q_B^+$ (e) to generate the background colors.

To display the overlaid vector field, however, requires a more involved formula. We use the exact same reactive current formula as in the supplement of [592021Strahan et al.Strahan, Antoszewski, Lorpaiboon, Vani, Weare,, and Dinner], but repeat it here for reference. The projected current is defined as

$$\mathbf{J}_{AB}^{\mathbf{Y}}(\mathbf{y}) = \int \mathbf{J}_{AB}(\mathbf{x}) \cdot \nabla\mathbf{Y}(\mathbf{x})\pi(\mathbf{x})\delta(\mathbf{Y}(\mathbf{x}) - \mathbf{y})\,d\mathbf{x} \tag{82}$$

In the discretized $\mathbf{y}$ space, this leads to the discretized projected current:

$$\mathbf{J}_{AB}^{\mathbf{Y}}(\mathbf{y}) \approx \frac{1}{2\Delta t}\sum_{n=1}^{N}\frac{d\pi}{d\mu}(\mathbf{X}_n(0))\Bigg[\mathbb{1}_{d\mathbf{y}}(\mathbf{X}_n(0))q_A^-(\mathbf{X}_n(0))q_B^+(\mathbf{X}_n(\theta_n))\frac{\mathbf{Y}(\mathbf{X}_n(\theta_n)) - \mathbf{Y}(\mathbf{X}_n(0))}{\theta_n} \tag{83}$$

$$+\mathbb{1}_{d\mathbf{y}}(\mathbf{X}_n(\Delta t))q_A^-(\mathbf{X}_n(\widetilde{\theta}_n))q_B^+(\mathbf{X}_n(\Delta t))\frac{\mathbf{Y}(\mathbf{X}_n(\Delta t)) - \mathbf{Y}(\mathbf{X}_n(\widetilde{\theta}_n))}{\Delta t - \widetilde{\theta}_n}\Bigg] \tag{84}$$

where $\theta_n$ and $\widetilde{\theta}_n$ are the "first-entry times" to $D = (A \cup B)^c$ in the $n$th trajectory with time running forward and backward, respectively. To visualize $\mathbf{J}_{AA}$, $\mathbf{J}_{BA}$, and $\mathbf{J}_{BB}$, we swap symbols accordingly on the committor subscripts. For the steady-state current $\mathbf{J}$, we replace all committors with 1.

We make use of the same formula to compute the flux distribution across surfaces of constant zonal wind ($\mathbf{Y} = U(30\,\text{km})$) in Fig. 6, or constant committor ($\mathbf{Y} = q_B^+$) in Figs. 7-10. The first and second terms in brackets represent contributions from the forward and backward generators, respectively. The flux contribution from the $n$th trajectory segment, therefore, is the average of those two terms. In selecting a set of candidate trajectories for a given level set, say of $q_B^+$, we choose trajectories whose *midpoint* has a $q_B^+$ value close to that level. This defines the discrete sampling of flux density, from which we construct medians and quantiles over various observables. In practice, we choose all points within a small range of $\mathbf{Y}$ values to aggregate statistics together; as the data size increases, this error window could be shrunk. Fig. 6 uses ranges of $U(30\,\text{km})$ equal to half the distance between successive levels. Figs. 7-10 uses ranges of $q_B^+$ equal to the distance between successive levels. In 10, the final committor level actually contains a range of $U(30\,\text{km})$, but we explicitly enforce the final boundary condition $\eta_B^+ = 0$ when $q_B^+ = 0$ at the endpoint of the plot, in order to make the quantile envelopes shrink to zero width.

## 2.5 Algorithmic parameters

Having sketched the general numerical procedure, we now provide the exact parameters used here, which are similar to those in [202021Finkel et al.Finkel, Webber, Abbot, Gerber,, and Weare]. We use $N = 3 \times 10^5$ trajectories, each of length $\Delta t = 20$ days, with a sampling interval of 1 day. The initial conditions are resampled from a long ($2 \times 10^5$-day) control simulation to be uniformly distributed on the space ($|\Psi|(30\,\text{km}), U(30\,\text{km})$). With a more complex, expensive model, we cannot rely on a control simulation to seed the initial points, but here we focus on TPT and DGA as a proof of concept rather than optimizing the numerical procedure. We use $M = 1500$ basis functions defined as indicators on a partition induced by $K$-means clustering on $\{\mathbf{X}_n(0)\}$. The clustering is hierarchical so that the cluster size does not become too imbalanced.

There are many potential directions for methodological improvement. In an expensive model without the ability to run a long control simulation, we should use a splitting and killing method to seed the initial trajectories across state space. Moreover, we could perform DGA repeatedly with new data seeded at each iteration in areas of high sensitivity. The choice of basis function can also powerfully affect DGA's performance. Indicator functions are advantageous in producing a bona fide Markov matrix and guaranteeing a maximum principle for the committor probabilities. However, smooth and/or global basis functions have in some cases found to be more efficient at capturing the structure of the committor with fewer basis elements [602019Thiede et al.Thiede, Giannakis, Dinner,, and Weare, 592021Strahan et al.Strahan, Antoszewski, Lorpaiboon, Vani, Weare,, and Dinn

# 3 Minimum-action method

To compute the minimum-action paths, we use a completely discrete approach for simplicity and to accommodate the low-rank nature of the stochastic forcing. Heuristically, we wish to find the most probable path connecting $A$ and $B$, which we take as the mode of the (discretized) path density over the distribution of paths from $A$ to $B$. For concreteness, fix $\mathbf{x}(0) = \mathbf{x}_0 \in A$ and a time horizon $T$ discretized into $K$ intervals, with a timestep $\delta t = T/K = 0.005$ days. The discretized dynamics evolve according to the Euler-Maruyama method as

$$\mathbf{x}(k\delta t) = \mathbf{x}((k-1)\delta t) + \boldsymbol{v}\big(\mathbf{x}((k-1)\Delta t)\big)\Delta t + \sigma\boldsymbol{\eta}_k\sqrt{\delta t} \tag{85}$$

where $\boldsymbol{\eta}_k$ is a vector of i.i.d. unit normal samples, $\boldsymbol{v}$ is the deterministic drift, and $\sigma \in \mathbb{R}^{d \times m}$ is the diffusion matrix, with a noise rank $m = 3$ and a spatially smooth structure as defined in [202021Finkel et al.Finkel, Webber, Abbot, Gerber,, and Weare]. In the classical minimum-action approach, $\sigma$ is assumed to be a $d \times d$ invertible matrix, and the probability density of a path $(\mathbf{x}_0, \ldots, \mathbf{x}_K)$ (where $\mathbf{x}_k = \mathbf{x}(k\delta t)$) is

$$\prod_{k=1}^{K} \mathcal{N}\Big(\mathbf{x}_k \Big| \mathbf{x}_{k-1} + \boldsymbol{v}(\mathbf{x}_{k-1})\delta t, \sigma\sigma^\top \delta t\Big) \tag{86}$$

$$= \prod_{k=1}^{K} \frac{1}{(2\pi\delta t)^{dK/2}(\det\sigma)^K} \times \tag{87}$$

$$\exp\left\{-\frac{1}{2}\Big(\mathbf{x}_k - \mathbf{x}_{k-1} - \boldsymbol{v}(\mathbf{x}_{k-1})\delta t\Big)^\top \frac{(\sigma\sigma^\top)^{-1}}{\delta t}\Big(\mathbf{x}_k - \mathbf{x}_{k-1} - \boldsymbol{v}(\mathbf{x}_{k-1})\delta t\Big)\right\} \tag{88}$$

$$\propto \exp\left\{-\frac{\delta t}{2}\sum_{k=1}^{K}\left(\frac{\mathbf{x}_k - \mathbf{x}_{k-1}}{\delta t} - \boldsymbol{v}(\mathbf{x}_{k-1})\right)^\top (\sigma\sigma^\top)^{-1}\left(\frac{\mathbf{x}_k - \mathbf{x}_{k-1}}{\delta t} - \boldsymbol{v}(\mathbf{x}_{k-1})\right)\right\} \tag{89}$$

$$\sim \exp\left\{-\frac{1}{2}\int_0^T \Big[\dot{\mathbf{x}}(t) - \boldsymbol{v}(\mathbf{x}(t))\Big](\sigma\sigma^\top)^{-1}\Big[\dot{\mathbf{x}}(t) - \boldsymbol{v}(\mathbf{x}(t))\Big]dt\right\} \text{ as } \delta t \to 0 \tag{90}$$

and the problem becomes to minimize the quadratic form in the argument of the exponential, which is the Freidlin-Wentzell action functional, subject to the constraint $\mathbf{x}_K \in B$. However, because we stir the wind field with smooth spatial forcing in only $m \ll d$ wavenumbers, $\sigma$ is low-rank and thus $\sigma\sigma^\top$ is singular. Given any realized path $(\mathbf{x}_0, \ldots, \mathbf{x}_K)$, there may be no possible underlying forcing $\boldsymbol{\eta}_k$ that could have produced it under our noise model. So the obvious optimization strategy of fixing $\mathbf{x}_0$ and $\mathbf{x}_K$ and varying the steps in between may lead to impossible paths. For this reason, we perform optimization in the space of perturbations, and ensure that every step of the optimization is realizable under our noise model. This is a strategy we adopt from the cyclogenesis model [**?**]. The result will be a simpler, convex objective function at the expense of a more complicated constraint. The probability density of a particular forcing sequence $(\boldsymbol{\eta}_1, \ldots, \boldsymbol{\eta}_K)$ is given by

$$\prod_{k=1}^{K} \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}\boldsymbol{\eta}_k^\top\boldsymbol{\eta}_k\right) = \frac{1}{(2\pi)^{mK/2}} \exp\left(-\frac{1}{2}\sum_{k=1}^{K}\boldsymbol{\eta}_k^\top\boldsymbol{\eta}_k\right) \tag{91}$$

The objective inside the exponential is now a simple quadratic in perturbation space which can be easily differentiated with respect to those perturbations. The constraint, meanwhile, takes the form of a complicated iterated function. Define the flow map $F(\mathbf{x}) = \mathbf{x} + \boldsymbol{v}(\mathbf{x})\delta t$ as the deterministic part of the timestep, so $\mathbf{x}_k = F(\mathbf{x}_{k-1}) + \sigma\boldsymbol{\eta}_k\sqrt{\delta t}$. In terms of $F$, the endpoint has to be written as a recursive function

$$\mathbf{x}_K = F(\mathbf{x}_{K-1}) + \sigma\boldsymbol{\eta}_K\sqrt{\delta t} \tag{92}$$

$$\mathbf{x}_{K-1} = F(\mathbf{x}_{K-2}) + \sigma\boldsymbol{\eta}_{K-1}\sqrt{\delta t} \tag{93}$$

$$\vdots \tag{94}$$

$$\mathbf{x}_1 = F(\mathbf{x}_0) + \sigma\boldsymbol{\eta}_1\sqrt{\delta t} \tag{95}$$

We impose the end constraint by adding to the action a penalty $\Phi(\mathbf{x}_K) = \text{dist}(\mathbf{x}_K, B)$, a function which linearly increases with distance to $B$. The full optimization problem is

$$\min_{\boldsymbol{\eta}}\left\{\frac{1}{2K}\sum_{k=1}^{K}\boldsymbol{\eta}_k^\top\boldsymbol{\eta}_k + \alpha\Phi(\mathbf{x}_K)\right\} \tag{96}$$

$$\mathbf{x}_0 \in A \text{ is fixed} \tag{97}$$

$$\mathbf{x}_k = F(\mathbf{x}_{k-1}) + \sigma\boldsymbol{\eta}_k\sqrt{\delta t} \text{ for } k = 1, \ldots, K \tag{98}$$

Here $\alpha$ is a weight which can be increased to harden the end constraint. We divide by $K$ so that the path action does not overwhelm the endpoint penalty as $K \to \infty$. (This makes the sum converge to an integral.) We set $\mathbf{x}_0$ to be the fixed point $\mathbf{a} \in A$ when finding the least-action path from $A$ to $B$ and the fixed point $\mathbf{b} \in B$ when finding the least-action path from $B$ to $A$. We used the L-BFGS method as implemented in `scipy`, with a maximum of 10 iterations. We differentiate $\Phi(x_K)$ with respect to $\boldsymbol{\eta}_k$ using knowledge of the adjoint model, with a backward pass through the path to compute each gradient. At each descent step, we refine the stepsize with backtracking line search. One way to guarantee the end constraint is ultimately satisfied is to gradually increase $\alpha$ and lengthen $T$; however, we found it sufficient to fix $\alpha = 1.0$ and $T = 100$, in keeping with the typical observed transit time. We have kept the algorithm simple, not devoting too much effort to finding the global optimimum over all time horizons, as we only care for a qualitative assessment to compare with results of TPT.

# References

[E et al., 2019] E, W., Li, T., and Vanden-Eijnden, E. (2019). *Applied stochastic analysis*, volume 199. American Mathematical Soc.

[E and Vanden-Eijnden, 2006] E, W. and Vanden-Eijnden, E. (2006). Towards a Theory of Transition Paths. *Journal of Statistical Physics*, 123(3):503.

[Finkel et al., 2021] Finkel, J., Webber, R. J., Abbot, D. S., Gerber, E. P., and Weare, J. (2021). Learning forecasts of rare stratospheric transitions from short simulations.

[Fitzsimmons and Pitman, 1999] Fitzsimmons, P. and Pitman, J. (1999). Kac's moment formula and the feynman–kac formula for additive functionals of a markov process. *Stochastic Processes and their Applications*, 79(1):117–134.

[Karatzas and Shreve, 1998] Karatzas, I. and Shreve, S. E. (1998). *Brownian Motion and Stochastic Calculus*. Springer.

[Metzner et al., 2006] Metzner, P., Schutte, C., and Vanden-Eijnden, E. (2006). Illustration of transition path theory on a collection of simple examples. *The Journal of Chemical Physics*, 125(8):1–2.

[Metzner et al., 2009] Metzner, P., Schutte, C., and Vanden-Eijnden, E. (2009). Transition path theory for markov jump processes. *Multiscale Modeling and Simulation*, 7(3):1192–1219.

[Oksendal, 2003] Oksendal, B. (2003). *Stochastic Differential Equations: An Introduction with Applications*. Springer.

[Plotkin et al., 2019] Plotkin, D. A., Webber, R. J., O'Neill, M. E., Weare, J., and Abbot, D. S. (2019). Maximizing simulated tropical cyclone intensity with action minimization. *Journal of Advances in Modeling Earth Systems*, 11(4):863–891.

[Strahan et al., 2021] Strahan, J., Antoszewski, A., Lorpaiboon, C., Vani, B. P., Weare, J., and Dinner, A. R. (2021). Long-time-scale predictions from short-trajectory data: A benchmark analysis of the trp-cage miniprotein. *Journal of Chemical Theory and Computation*, 17(5):2948–2963. PMID: 33908762.

[Thiede et al., 2019] Thiede, E., Giannakis, D., Dinner, A. R., and Weare, J. (2019). Approximation of dynamical quantities using trajectory data. *arXiv:1810.01841 [physics.data-an]*, pages 1–24.

[Vanden-Eijnden, 2006] Vanden-Eijnden, E. (2006). *Transition Path Theory*, pages 453–493. Springer Berlin Heidelberg, Berlin, Heidelberg.

[Vanden-Eijnden and E, 2010] Vanden-Eijnden, E. and E, W. (2010). Transition-path theory and path-finding algorithms for the study of rare events. *Annual Review of Physical Chemistry*, 61(1):391–420.