

1 **Learning forecasts of rare stratospheric transitions from short simulations**

2 Justin Finkel*

3 *Committee on Computational and Applied Mathematics, University of Chicago*

4 Robert J. Webber

5 *Courant Institute of Mathematical Sciences, New York University*

6 Edwin P. Gerber

7 *Courant Institute of Mathematical Sciences, New York University*

8 Dorian S. Abbot

9 *Department of the Geophysical Sciences, University of Chicago*

10 Jonathan Weare

11 *Courant Institute of Mathematical Sciences, New York University*

12 *Corresponding author: Justin Finkel, jfinkel@uchicago.edu

ABSTRACT

13 Rare events arising in nonlinear atmospheric dynamics remain hard to predict and attribute. We
14 address the problem of forecasting rare events in a prototypical example, Sudden Stratospheric
15 Warmings (SSWs). Approximately once every other winter, the boreal stratospheric polar vortex
16 rapidly breaks down, shifting midlatitude surface weather patterns for months. We focus on two
17 key quantities of interest: the probability of an SSW occurring, and the expected lead time if it does
18 occur, as functions of initial condition. These *optimal forecasts* concretely measure the event's
19 progress. Direct numerical simulation can estimate them in principle, but is prohibitively expensive
20 in practice: each rare event requires a long integration to observe, and the cost of each integration
21 grows with model complexity. We describe an alternative approach using integrations that are
22 *short* compared to the timescale of the warming event. We compute the probability and lead time
23 efficiently by solving equations involving the transition operator, which encodes all information
24 about the dynamics. We relate these optimal forecasts to a small number of interpretable physical
25 variables, suggesting optimal measurements for forecasting. We illustrate the methodology on a
26 prototype SSW model developed by Holton and Mass (1976) and modified by stochastic forcing.
27 While highly idealized, this model captures the essential nonlinear dynamics of SSWs and exhibits
28 the key forecasting challenge: the dramatic separation in timescales between a single event and
29 the return time between successive events. Our methodology is designed to fully exploit high-
30 dimensional data from models and observations, and can identify detailed predictors of many
31 complex rare events in meteorology.

32 **1. Introduction**

33 As computing power increases and weather models grow more intricate and capable of generating
34 a vast wealth of realistic data, the goal of extreme weather event prediction appears less distant
35 (Vitart and Robertson 2018). To take full advantage of the increased computing power, we must
36 develop new approaches to efficiently manage and parse the data we generate (or observe) to
37 derive physically interpretable, actionable insights. Extreme weather events are worthy targets
38 for simulation owing to their destructive potential to life and property. Rare events have attracted
39 significant simulation efforts recently, including hurricanes (e.g., Zhang and Sippel 2009; Webber
40 et al. 2019; Plotkin et al. 2019), heat waves (e.g., Ragone et al. 2018), rogue waves (e.g., Dematteis
41 et al. 2018), and space weather events (e.g., coronal mass ejections; Ngwira et al. (2013)). These
42 are very difficult to characterize and predict, being exceptionally rare and pathological outliers
43 in the spectrum of weather events. Ensemble forecasting in numerical weather prediction is best
44 suited to estimate statistics of the average or most likely scenarios, and specialized methods are
45 needed to examine the more extreme outlier scenarios.

46 In this study, we advance an alternative computational approach to predicting and understanding
47 general rare events without sacrificing model fidelity. Our method relies on data generated by a
48 high-fidelity model with a state space with many degrees of freedom d , representing, for example,
49 spatial resolution of the primitive equations. In this way, our method is similar to recently introduced
50 reduced order modeling techniques using statistical and machine learning (e.g., Kashinath et al.
51 (2021) and references therein). However, in contrast to other data-driven techniques, our approach
52 focuses on directly computing key quantities of interest that characterize the essential predictability
53 of the rare event, rather than trying to capture the full detailed evolution of the system. In particular,
54 we will compute estimators of *statistically optimal forecasts* that are useful for initial conditions

55 somewhere between a “typical” configuration A and an “anomalous” configuration B that defines
56 the rare event, where typical and anomalous are user-defined. We focus on two forecasts in
57 particular to quantify risk. The *committor* is the probability that a given initial condition evolves
58 directly into B rather than A . Given that it does reach B first, the *conditional mean first passage*
59 *time*, or *lead time*, is the expected time that it takes to get there. The committor appears prominently
60 in the molecular dynamics literature, with some recent applications in geoscience including Tantet
61 et al. (2015); Lucente et al. (2019), and Finkel et al. (2020), which compute the committor for
62 low-dimensional atmospheric models.

63 Both quantities depend on the initial condition, defining functions over d -dimensional state space
64 that encode important information regarding the fundamental causes and precursors of the rare
65 event. However, “decoding” the physical insights is not automatic. With real-time measurement
66 constraints, the risk metrics must be estimated from low-dimensional proxies. Even visualizing
67 them requires projecting down to one or two dimensions. This calls for a principled selection of
68 low-dimensional coordinates which are both physically meaningful and statistically informative
69 for our chosen risk metrics. We address this problem using sparse regression, a simple but easily
70 extensible solution with the potential to inform optimal measurement strategies to estimate risk as
71 precisely as possible under constraints.

72 Estimation of the committor and lead time is a challenge. We employ a method that uses
73 a large data set of short-time independent simulations. We represent the committor and lead
74 time as solutions to Feynman-Kac formulae (Oksendal 2003), which relate long-time forecasts to
75 instantaneous tendencies. These equations are elegant and general, but computationally daunting:
76 in the continuous time and space limit, they become partial differential equations (PDE) with d
77 independent variables—the same as the model state space dimension. It is therefore hopeless to
78 solve the equations using any standard spatial discretization. But, as we demonstrate, the equations

79 can be solved with remarkable accuracy by expanding in a basis of functions informed by the data
80 set.

81 We illustrate our approach on the highly simplified Holton-Mass model (Holton and Mass
82 1976; Christiansen 2000) with stochastic velocity perturbations in the spirit of Birner and Williams
83 (2008). The Holton-Mass model is well-understood dynamically in light of decades of analysis and
84 experiments, yet complex enough to present the essential computational difficulties of probabilistic
85 forecasting and test our methods for addressing them. In particular, this system captures the
86 key difficulty in sampling rare events. The vast majority of the time, the system sits in one of
87 two metastable states, characterizing a strong or weak vortex respectively. Extreme events are
88 the infrequent jumps from one state to another. Our computational framework can accurately
89 characterize these rare transitions using only a data set of “short” model simulations, short not
90 only compared to the long periods the system sits in one state or the other, but also relative to
91 the timescale of the transition events themselves. In the future, the same methodology could be
92 applied to query the properties of more complex models, such as GCMs, where less theoretical
93 understanding is available.

94 In section 2, we review the dynamical model and define the specific rare event of interest. In
95 section 3, we formally define the risk metrics introduced above and visualize the results for the
96 Holton-Mass model, including a discussion of physical and practical insights gleaned from our
97 approach. In section 4 we identify an optimal set of reduced coordinates for estimating risk using
98 sparse regression. These results will provide motivation for the computational method, which we
99 present afterward in section 5 along with accuracy tests. We then lay out future prospects and
100 conclude in section 6.

101 2. Holton-Mass model

102 Holton and Mass (1976) devised a simple model of the stratosphere aimed at reproducing
103 observed intra-seasonal oscillations of the polar vortex, which they termed “stratospheric vacil-
104 lation cycles.” Earlier SSW models, originating with that of Matsuno (1971), proposed upward-
105 propagating planetary waves as the major source of disturbance to the vortex. While Matsuno
106 (1971) used impulsive forcing from the troposphere as the source of planetary waves, Holton
107 and Mass (1976) suggested that even stationary tropospheric forcing could lead to an oscillatory
108 response, suggesting that the stratosphere can self-sustain its own oscillations. While the Holton-
109 Mass model is meant to represent internal stratospheric dynamics, Sjoberg and Birner (2014) point
110 out that the stationary boundary condition does not lead to stationary wave activity flux, meaning
111 that even the Holton-Mass model involves some dynamic interaction between the troposphere and
112 stratosphere. Isolating internal from external dynamics is a subtle modeling question, but in the
113 present paper we adhere to the original Holton-Mass framework for simplicity. Our methodology
114 applies equally well to other formulations.

115 Radiative cooling through the stratosphere and wave perturbations at the tropopause are the two
116 competing forces that drive the vortex in the Holton-Mass model. Altitude-dependent cooling
117 relaxes the zonal wind toward a strong vortex in thermal wind balance with a radiative equilibrium
118 temperature field. Gradients in potential vorticity along the vortex, however, can allow the propaga-
119 tion of Rossby waves. When conditions are just right, a Rossby wave emerges from the tropopause
120 and rapidly propagates upward, sweeping heat poleward and stalling the vortex by depositing a
121 burst of negative momentum. The vortex is destroyed and begins anew the rebuilding process.

122 Yoden (1987a) found that for a certain range of parameter settings, these two effects balance each
123 other to create two distinct stable regimes: a strong vortex with zonal wind close to the radiative

124 equilibrium profile, and a weak vortex with a possibly oscillatory wind profile. We focus our study
 125 on this bistable setting as a prototypical model of atmospheric regime behavior. The transition
 126 from strong to weak vortex state captures the essential dynamics of an SSW.

127 The Holton-Mass model takes the linearized quasigeostrophic potential vorticity (QGPV) equa-
 128 tion for a perturbation streamfunction $\psi'(x, y, z, t)$ on top of a zonal mean flow $\bar{u}(y, z, t)$, and
 129 projects these two fields onto a single zonal wavenumber $k = 2/(a \cos 60^\circ)$ and a single meridional
 130 wavenumber $\ell = 3/a$, where a is the Earth's radius. This notation is consistent with Holton and
 131 Mass (1976) and Christiansen (2000), and we refer the reader to these earlier papers for complete
 132 description of the equations and parameters. The resulting ansatz is

$$\bar{u}(y, z, t) = U(z, t) \sin(\ell y) \quad (1)$$

$$\psi'(x, y, z, t) = \text{Re}\{\Psi(z, t)e^{ikx}\}e^{z/2H} \sin(\ell y)$$

133 which is fully determined by the reduced state space $U(z, t)$, and $\Psi(z, t)$, the latter being complex-
 134 valued. H is a scale height, 7 km. Inserting this into the linearized QGPV equations yields the
 135 coupled PDE system

$$\begin{aligned} & \left[-\left(\mathcal{G}^2(k^2 + \ell^2) + \frac{1}{4}\right) + \frac{\partial^2}{\partial z^2} \right] \frac{\partial \Psi}{\partial t} \\ & = \left[\left(\frac{\alpha}{4} - \frac{\alpha_z}{2} - i\mathcal{G}^2 k \beta \right) - \alpha_z \frac{\partial}{\partial z} - \alpha \frac{\partial^2}{\partial z^2} \right] \Psi \\ & + \left\{ ik\varepsilon \left[\left(k^2 \mathcal{G}^2 + \frac{1}{4} \right) - \frac{\partial}{\partial z} + \frac{\partial^2}{\partial z^2} \right] U \right\} \Psi - ik\varepsilon \frac{\partial^2 \Psi}{\partial z^2} U \end{aligned} \quad (2)$$

136 for $\Psi(z, t)$, and

$$\begin{aligned} & \left(-\mathcal{G}^2 \ell^2 - \frac{\partial}{\partial z} + \frac{\partial^2}{\partial z^2} \right) \frac{\partial U}{\partial t} = [(\alpha_z - \alpha)U_z^R - \alpha U_{zz}^R] \\ & - \left[(\alpha_z - \alpha) \frac{\partial}{\partial z} + \alpha \frac{\partial^2}{\partial z^2} \right] U + \frac{\varepsilon k \ell^2}{2} e^z \text{Im} \left\{ \Psi \frac{\partial^2 \Psi^*}{\partial z^2} \right\} \end{aligned} \quad (3)$$

137 for $U(z, t)$. Here, $\varepsilon = 8/(3\pi)$ is a coefficient for projecting $\sin^2(\ell y)$ onto $\sin(\ell y)$. We have
 138 nondimensionalized the equations with the parameter $\mathcal{G}^2 = H^2 N^2 / (f_0^2 L^2)$, where $N^2 = 4 \times 10^{-4} \text{ s}^{-2}$

139 is a constant stratification (Brunt-Väisälä frequency), f_0 is the Coriolis parameter, and $L = 2.5 \times 10^5$
 140 m is a horizontal length scale, selected in order to create a homogeneously shaped data set more
 141 suited to our analysis. See Holton and Mass (1976); Yoden (1987a); Christiansen (2000) for details
 142 on parameters. Boundary conditions are prescribed at the bottom of the stratosphere, which in this
 143 model corresponds to $z = 0$ km, and the top of the stratosphere $z_{top} = 70$ km.

$$\begin{aligned} \Psi(0, t) &= \frac{gh}{f_0}, & \Psi(z_{top}, t) &= 0, \\ U(0, t) &= U^R(0), & \partial_z U(z_{top}, t) &= \partial_z U^R(z_{top}). \end{aligned} \quad (4)$$

144 The vortex-stabilizing influence is represented by $\alpha(z)$, the altitude-dependent cooling coefficient,
 145 and the radiative wind profile $U^R(z) = U^R(0) + \frac{\gamma}{1000}z$ (with z in m), which relaxes the vortex toward
 146 radiative equilibrium. Here $\gamma = \mathcal{O}(1)$ is the vertical wind shear in m/s/km. The competing force
 147 of wave perturbation is encoded through the lower boundary condition $\Psi(0, t) = gh/f_0$.

148 Detailed bifurcation analysis of the model by both Yoden (1987a) and Christiansen (2000) in
 149 (γ, h) space revealed the bifurcations that lead to bistability, vacillations, and ultimately quasiperi-
 150 odicity and chaos. Here we will focus on an intermediate parameter setting of $\gamma = 1.5$ m/s/km
 151 and $h = 38.5$ m, where two stable states coexist: a strong vortex with U closely following U^R and
 152 an almost barotropic stationary wave, as well as a weak vortex with U dipping close to zero at
 153 an intermediate altitude and a stationary wave with strong westward phase tilt. The two stable
 154 equilibria, which we call **a** and **b**, are illustrated in Figure 1(a,b) by their z -dependent zonal wind
 155 and perturbation streamfunction profiles.

156 The two equilibria can be interpreted as two different winter climatologies, one with a strong
 157 vortex and one with a weak vortex susceptible to vacillation cycles. To explore transitions between
 158 these two states, we follow Birner and Williams (2008) and modify the Holton-Mass equations
 159 with small additive noise in the U variable to mimic momentum perturbations by smaller scale

160 Rossby waves, gravity waves, and other unresolved sources. The form of noise will be specified in
 161 Equation (7).

162 While the details of the additive noise are ad hoc, the general approach can be more rigorously
 163 justified through the Mori-Zwanzig formalism (Zwanzig 2001). Because many hidden degrees
 164 of freedom are being projected onto the low-dimensional space of the Holton-Mass model, the
 165 dynamics on small observable subspaces can be considered stochastic. This is the perspective
 166 taken in stochastic parameterization of turbulence and other high-dimensional chaotic systems
 167 (Hasselmann 1976; DelSole and Farrell 1995; Franzke and Majda 2006; Majda et al. 2001; Gottwald
 168 et al. 2016). In general, unobserved deterministic dynamics can make the system non-Markovian,
 169 which technically violates the assumptions of our methodology. However, with sufficient separation
 170 of timescales the Markovian assumption is not unreasonable. Furthermore, memory terms can
 171 be ameliorated by lifting data back to higher-dimensional state space with time-delay embedding
 172 (Berry et al. 2013; Thiede et al. 2019; Lin and Lu 2021).

173 We follow Holton and Mass (1976) and discretize the equations using a finite-difference method
 174 in z , with 27 vertical levels (including boundaries). After constraining the boundaries, there are
 175 $d = 3 \times (27 - 2) = 75$ degrees of freedom in the model. Christiansen (2000) investigated higher
 176 resolution and found negligible differences. The full discretized state is represented by a long
 177 vector

$$\begin{aligned} \mathbf{X}(t) = & \left[\operatorname{Re}\{\Psi\}(\Delta z, t), \dots, \operatorname{Re}\{\Psi\}(z_{top} - \Delta z, t), \right. \\ & \operatorname{Im}\{\Psi\}(\Delta z, t), \dots, \operatorname{Im}\{\Psi\}(z_{top} - \Delta z, t), \\ & \left. U(\Delta z, t), \dots, U(z_{top} - \Delta z, t) \right] \in \mathbb{R}^d = \mathbb{R}^{75} \end{aligned} \quad (5)$$

178 The deterministic system can be written $d\mathbf{X}(t)/dt = \mathbf{v}(\mathbf{X}(t))$ for a vector field $\mathbf{v} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ specified
 179 by discretizing (2) and (3). Under deterministic dynamics, $\mathbf{X}(t) \rightarrow \mathbf{a}$ or $\mathbf{X}(t) \rightarrow \mathbf{b}$ as $t \rightarrow \infty$

180 depending on initial conditions. The addition of white noise changes the system into an Itô
 181 diffusion

$$d\mathbf{X}(t) = \mathbf{v}(\mathbf{X}(t)) dt + \boldsymbol{\sigma}(\mathbf{X}(t)) d\mathbf{W}(t) \quad (6)$$

182 where $\boldsymbol{\sigma} : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$ imparts a correlation structure to the vector $\mathbf{W}(t) \in \mathbb{R}^m$ of independent
 183 standard white noise processes. As discussed above, we design $\boldsymbol{\sigma}$ to be a low-rank, constant matrix
 184 that adds spatially smooth stirring to only the zonal wind U (not the streamfunction Ψ) and which
 185 respects boundary conditions at the bottom and top of the stratosphere. Its structure is defined
 186 by the following Euler-Maruyama scheme: in a timesetep $\delta t = 0.005$ days, after a deterministic
 187 forward Euler step we add the stochastic perturbation to zonal wind on large vertical scales

$$\delta U(z) = \sigma_U \sum_{k=0}^m \eta_k \sin \left[\left(k + \frac{1}{2} \right) \pi \frac{z}{z_{top}} \right] \sqrt{\delta t} \quad (7)$$

188 where η_k ($k = 0, 1, 2$) are independent unit normal samples, $m = 2$, and σ_U is a scalar that sets the
 189 magnitudes of entries in $\boldsymbol{\sigma}$. In terms of physical units,

$$\sigma_U^2 = \frac{\mathbb{E}[(\delta U)^2]}{\delta t} \approx (1 \text{ m/s})^2 / \text{day} \quad (8)$$

190 σ_U has units of $(L/T)/T^{1/2}$, where the square-root of time comes from the quadratic variation of the
 191 Wiener process. It is best interpreted in terms of the daily root-mean-square velocity perturbation
 192 of 1.0 m/s. We have experimented with this value, and found that reducing the noise level below
 193 0.8 dramatically reduces the frequency of transitions, while increasing it past 1.5 washes out
 194 metastability. We keep σ_U constant going forward as a favorable numerical regime to demonstrate
 195 our approach, while acknowledging that the specifics of stochastic parameterization are important
 196 in general to obtain accurate forecasts. The resulting matrix $\boldsymbol{\sigma}$ is 75×3 , with nonzero entries only
 197 in the last 25 rows as forcing only applies to $U(z)$.

198 A long simulation of the model reveals metastability, with the system tending to remain close to
 199 one fixed point for a long time before switching quickly to the other, as shown by the time series

200 of $U(30\text{km})$ in panel (d) of Figure 1. Panel (e) shows a projection of the steady state distribution,
 201 also known as the equilibrium/invariant distribution, of U as a function of z . We call this density
 202 $\pi(\mathbf{x})$, which is a function over the full d -dimensional state space. We focus on the zonal wind
 203 U at 30 km following Christiansen (2000), because this is where its strength is minimized in the
 204 weak vortex. While the two regimes are clearly associated with the two fixed points, they are better
 205 characterized by extended *regions* of state space with strong and weak vortices. We thus define the
 206 two metastable subsets of \mathbb{R}^d

$$A = \{\mathbf{X} : U(\mathbf{X})(30\text{km}) \geq U(\mathbf{a})(30\text{km}) = 53.8\text{ m/s}\},$$

$$B = \{\mathbf{X} : U(\mathbf{X})(30\text{km}) \leq U(\mathbf{b})(30\text{km}) = 1.75\text{ m/s}\}.$$

207 This straightforward definition roughly follows the convention of Charlton and Polvani (2007),
 208 which defines an SSW as a reversal of zonal winds at 10 hPa. We use 30 km for consistency with
 209 Christiansen (2000); this is technically higher than 10 hPa because $z = 0$ in the Holton-Mass model
 210 represents the tropopause. Our method is equally applicable to any definition, and the results
 211 are not qualitatively dependent on this choice. Incidentally, the analysis tools we present may be
 212 helpful in distinguishing predictability properties between different definitions. In fact, we will
 213 show that the height neighborhood of 20 km is actually more salient for predicting the event than
 214 wind at the 30-km level, even when the event is defined by wind at 30 km! This emerges from
 215 statistical analysis alone, and gives us confidence that essential SSW properties are stable with
 216 respect to reasonable changes in definition.

217 The orange highlights in Figure 1 (d) begin when the system exits the A region bound for B ,
 218 and end when the system enters B . The green highlights start when the system leaves B bound
 219 for A , and end when A is reached. Note that $A \rightarrow B$ transitions, SSWs, are much shorter in
 220 duration than $B \rightarrow A$ transitions. Figure 1 (c) shows the same paths, but viewed parametrically

221 in a two-dimensional state space consisting of integrated heat flux or IHF $\int_{0 \text{ km}}^{30 \text{ km}} e^{-z/H} \overline{v'T'} dz$, and
 222 zonal wind $U(30 \text{ km})$. IHF is an informative number because it captures both magnitude and phase
 223 information of the streamfunction in the Holton-Mass model:

$$\text{IHF} = \int_{0 \text{ km}}^{30 \text{ km}} e^{-z/H} \overline{v'T'} dz \propto \int_{0 \text{ km}}^{30 \text{ km}} |\Psi|^2 \frac{\partial \varphi}{\partial z} dz \quad (9)$$

224 where φ is the phase of Ψ . The $A \rightarrow B$ and $B \rightarrow A$ transitions are again highlighted in orange
 225 and green respectively, showing geometrical differences between the two directions. We will
 226 refer to the $A \rightarrow B$ transition as an SSW event, even though it is more accurately a transition
 227 between climatologies according to the Holton-Mass interpretation. The $B \rightarrow A$ transition is a
 228 vortex restoration event. Our focus in this paper is on predicting these transition events (mainly the
 229 $A \rightarrow B$ direction) and monitoring their progress in a principled way. In the next section we explain
 230 the formalism for doing so.

231 3. Forecast functions: the committor and lead time statistics

232 a. Defining risk and lead time

233 We will introduce the quantities of interest by way of example. First, suppose the stratosphere is
 234 observed in an initial state $\mathbf{X}(0) = \mathbf{x}$ that is neither in A nor B , so $U(\mathbf{b})(30 \text{ km}) < U(\mathbf{x})(30 \text{ km}) <$
 235 $U(\mathbf{a})(30 \text{ km})$ and the vortex is somewhat weakened, but not completely broken down. We call this
 236 intermediate zone $D = (A \cup B)^c$ (the complement of the two metastable sets). Because A and B
 237 are attractive, the system will soon find its way to one or the other at the *first-exit time* from D ,
 238 denoted

$$\tau_{D^c} = \min\{t \geq 0 : \mathbf{X}(t) \in D^c\} \quad (10)$$

239 Here, D^c emphasizes that the process has left D , i.e., gone to A or B . The first-exit location
 240 $\mathbf{X}(\tau_{D^c})$ is itself a random variable which importantly determines how the system exits D : either

241 $\mathbf{X}(\tau_{D^c}) \in A$, meaning the vortex restores to radiative equilibrium, or $\mathbf{X}(\tau_{D^c}) \in B$, meaning the
 242 vortex breaks down into vacillation cycles. A fundamental goal of forecasting is to determine the
 243 probabilities of these two events, which naturally leads to the definition of the (forward) committor
 244 function

$$q^+(\mathbf{x}) = \begin{cases} \mathbb{P}_{\mathbf{x}}\{\mathbf{X}(\tau_{D^c}) \in B\} & \mathbf{x} \in D = (A \cup B)^c \\ 0 & \mathbf{x} \in A \\ 1 & \mathbf{x} \in B \end{cases} \quad (11)$$

245 where the subscript \mathbf{x} indicates that the probability is conditional on a fixed initial condition
 246 $\mathbf{X}(0) = \mathbf{x}$, i.e., $\mathbb{P}_{\mathbf{x}}\{\cdot\} = \mathbb{P}\{\cdot | \mathbf{X}(0) = \mathbf{x}\}$. The superscript “+” distinguishes the forward committor
 247 from the *backward committor*, an analogous quantity for the time-reversed process which we do
 248 not use in this paper. Throughout, we will use capital $\mathbf{X}(t)$ to denote a stochastic process, and
 249 lower-case \mathbf{x} to represent a specific point in state space, typically an initial condition, i.e., $\mathbf{X}(0) = \mathbf{x}$.
 250 Both are $d = 75$ -dimensional vectors.

251 The committor is the probability that the system will be in state B (the disturbed state) next rather
 252 than A (the strong vortex state). Hence $q^+(\mathbf{x}) = 0$ if you start in A , and is 1 if you are already in B .
 253 In between (i.e., when $\mathbf{x} \in D$), $q^+(\mathbf{x})$ tells you the probability that you will first go to B rather than
 254 to A . That is, $q^+(\mathbf{x})$ tells you the probability that an SSW will happen.

255 Another important forecasting quantity is the lead time to the event of interest. While the forward
 256 committor reveals the probability of experiencing vortex breakdown *before* returning to a strong
 257 vortex, it does not say how long either event will take. Furthermore, even if the vortex is restored
 258 first, how long will it be until the next SSW does occur? The time until the next SSW event is
 259 denoted τ_B , again a random variable, whose distribution depends on the initial condition \mathbf{x} . We
 260 call $\mathbb{E}_{\mathbf{x}}[\tau_B]$ the *mean first passage time* (MFPT) to B . Conversely, we may ask how long a vortex

261 disturbance will persist before normal conditions return; the answer (on average) is $\mathbb{E}_{\mathbf{x}}[\tau_A]$, the
 262 mean first passage time to A . These same quantities have been calculated previously in other
 263 simplified models, e.g. Birner and Williams (2008) and Esler and Mester (2019).

264 $\mathbb{E}_{\mathbf{x}}[\tau_B]$ has an obvious shortcoming: it is an average over all paths starting from \mathbf{x} , including those
 265 which go straight into B (i.e., an orange trajectory in Figure 1c,d) and the rest which return to A i.e.,
 266 a green trajectory) and linger there, potentially for a very long time, before eventually re-crossing
 267 back into B . It is more relevant for near-term forecasting to condition τ_B on the event that an SSW
 268 is coming before the strong vortex returns. For this purpose, we introduce the *conditional* mean
 269 first passage time, or lead time, to B :

$$\eta^+(\mathbf{x}) := \mathbb{E}_{\mathbf{x}}[\tau_B | \tau_B < \tau_A] \quad (12)$$

270 which quantifies the suddenness of SSW.

271 All of these quantities can, in principle, be estimated by direct numerical simulation, or *shooting*.
 272 For example, suppose we observe an initial condition $\mathbf{X}(0) = \mathbf{x}$ in an operational forecasting
 273 setting, and wish to estimate the probability and lead time for the event of next hitting B . We would
 274 initialize an ensemble $\{\mathbf{X}_n(0) = \mathbf{x}, n = 1, \dots, N\}$ and evolve each member forward in time until it
 275 hits A or B at the random time τ_n . In an explicitly stochastic model, random forcing would drive
 276 each member to a different fate, while in a deterministic model their initial conditions would be
 277 perturbed slightly. To estimate the committor to B , we could calculate the fraction of members
 278 that hit B first. Averaging the arrival times (τ_n), over only those members gives an estimate of the
 279 lead time to B . For a single initial condition \mathbf{x} reasonably close to B , this direct shooting method
 280 may be the most economical. But how do we systematically compute $q^+(\mathbf{x})$ over all of state space
 281 (here 75 variables, but potentially billions of variables in a GCM or other state-of-the-art forecast
 282 system)?

283 For this more ambitious goal, the direct shooting method is prohibitively expensive. By definition,
284 transitions between A and B are infrequent. Therefore, if starting from \mathbf{x} far from B , a huge number
285 of sampled trajectories (N) will be required to observe even a small number ending in B , and they
286 may take a long time to get there. If instead we could precompute these functions offline over all
287 of state space, the online forecasting problem would reduce to “reading off” the committor and
288 lead time with every new observation. Achieving this goal is the key point of our paper, and we
289 achieve this using the dynamical Galerkin approximation, or DGA, recipe described by Thiede
290 et al. (2019).

291 A brute force way to estimate these functions is to integrate the model for a long time until
292 it reaches statistical steady state, meaning it has explored its attractor thoroughly according to
293 the steady state distribution. After long enough, it will have wandered close to every point \mathbf{x}
294 sufficiently often to estimate $q^+(\mathbf{x})$ and η^+ robustly as in shooting. We have performed such a
295 “control simulation” of 5×10^5 days for validation purposes, but our main contribution in this
296 paper is to compute the forecast functions using only *short* trajectories with DGA, allowing for
297 massive parallelization. However, we will defer the methodological details to Section 5, and first
298 justify the effort with some results. We visualize the committor and lead time computed from short
299 trajectories and elaborate on their interpretation, utility, and relationship to ensemble forecasting
300 methods.

301 *b. Steady state distribution*

302 Before visualizing the committor and lead time, it will be helpful to have a precise notion of the
303 steady state distribution, denoted $\pi(\mathbf{x})$, a probability density that describes the long-term behavior
304 of a stochastic process $\mathbf{X}(t)$. Assuming the system is ergodic, averages over time are equivalent to

305 averages over state space with respect to π . That is, for any well-behaved function $g : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T g(\mathbf{X}(t)) dt = \int_{\mathbb{R}^d} g(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x} =: \langle g \rangle_\pi \quad (13)$$

306 For example, if $g(\mathbf{x}) = \mathbb{1}_S(\mathbf{x})$ (an indicator function, which is 1 for $\mathbf{x} \in S \subset \mathbb{R}^d$ and 0 for $\mathbf{x} \notin S$),
 307 Equation (13) says that the fraction of time spent in S can be found by integrating the density
 308 over S . The density peaks in Figure 1(d) indicates clearly that the neighborhoods of \mathbf{a} and \mathbf{b}
 309 are two such regions with especially large probability under π . Note that both sides of (13) are
 310 independent of the initial condition, which is forgotten eventually. Short-term forecasts are by
 311 definition out-of-equilibrium processes, depending critically on initial conditions; however, $\pi(\mathbf{x})$
 312 is important to us here as a “default” distribution for missing information. If the initial condition is
 313 only partially observed, e.g. in only one coordinate, we have no information about the other $d - 1$
 314 dimensions, and in many cases the most principled tactic is to assume those other dimensions are
 315 distributed according to π , conditional on the observation.

316 *c. Visualizing committor and lead times*

317 The forecasts $q^+(\mathbf{x})$ and $\eta^+(\mathbf{x})$ are functions of a high-dimensional space \mathbb{R}^d . However, these
 318 degrees of freedom may not all be “observable” in a practical sense, given the sparsity and resolution
 319 limits of weather sensors, and visualizing them requires projecting onto reduced-coordinate spaces
 320 of dimension 1 or 2. We call these “collective variables” (CVs) following chemistry literature
 321 (e.g., Noé and Clementi 2017), and denote them as vector-valued functions from the full state
 322 space to a reduced space, $\theta : \mathbb{R}^d \rightarrow \mathbb{R}^k$, where $k = 1$ or 2 . For instance, Figure 1 (c) plots
 323 trajectories in the CV space consisting of integrated heat flux and zonal wind at 30 km: $\theta(\mathbf{x}) =$
 324 $\left(\int_0^{30\text{km}} e^{-z/H} \overline{v'T'} dz, U(30\text{km}) \right)$. The first component is a nonlinear function involving products
 325 of $\text{Re}\{\Psi\}$ and $\text{Im}\{\Psi\}$, while the second component is a linear function involving only U at a certain

326 altitude. For visualization in general, we have to approximate a function $F : \mathbb{R}^d \rightarrow \mathbb{R}$, such as the
 327 committor or lead time, as a function of reduced coordinates. That is, we wish to find $f : \mathbb{R}^k \rightarrow \mathbb{R}$
 328 such that $F(\mathbf{x}) \approx f(\boldsymbol{\theta}(\mathbf{x}))$. Given a fixed CV space $\boldsymbol{\theta}$, an “optimal” f is chosen by minimizing
 329 some function-space metric between $f \circ \boldsymbol{\theta}$ and F .

330 A natural choice is the mean-squared error weighted by the steady state distribution π , so the
 331 projection problem is to minimize over functions $f : \mathbb{R}^k \rightarrow \mathbb{R}$ the penalty

$$\begin{aligned} S[f; \boldsymbol{\theta}] &:= \|f \circ \boldsymbol{\theta} - F\|_{L^2(\pi)}^2 \\ &= \int_{\mathbb{R}^d} \left[f(\boldsymbol{\theta}(\mathbf{x})) - F(\mathbf{x}) \right]^2 \pi(\mathbf{x}) d\mathbf{x}. \end{aligned} \quad (14)$$

332 The optimal f for this purpose is the conditional expectation

$$\begin{aligned} f(\mathbf{y}) &= \mathbb{E}_{\mathbf{X} \sim \pi} [F(\mathbf{X}) | \boldsymbol{\theta}(\mathbf{X}) = \mathbf{y}] \\ &= \lim_{|d\mathbf{y}| \rightarrow 0} \frac{\int f(\mathbf{x}) \mathbb{1}_{d\mathbf{y}}(\boldsymbol{\theta}(\mathbf{x})) \pi(\mathbf{x}) d\mathbf{x}}{\int \mathbb{1}_{d\mathbf{y}}(\boldsymbol{\theta}(\mathbf{x})) \pi(\mathbf{x}) d\mathbf{x}} \end{aligned} \quad (15)$$

333 where $d\mathbf{y}$ is a small neighborhood about \mathbf{y} in CV space \mathbb{R}^k . The subscript $\mathbf{X} \sim \pi$ means that the
 334 expectation is with respect to a random variable \mathbf{X} distributed according to $\pi(\mathbf{x})$, i.e., at steady
 335 state. Figure 2 uses this formula to display one-dimensional projections of the committor (first row)
 336 and lead time (second row), as well as the one-standard deviation envelope incurred by projecting
 337 out the other 74 degrees of freedom. This “projection error” is defined as the square root of the
 338 conditional variance

$$V[f](\mathbf{y}) = \mathbb{E}_{\mathbf{X} \sim \pi} \left[\left(F(\mathbf{X}) - f(\mathbf{y}) \right)^2 \middle| \boldsymbol{\theta}(\mathbf{X}) = \mathbf{y} \right]. \quad (16)$$

339 Each quantity is projected onto two different one-dimensional CVs: U (30 km) (first column) and
 340 IHF (second column). In panel (a), for example, we see the committor is a decreasing function
 341 of U : the weaker the wind, the more likely a vortex breakdown. Moreover, the curve provides a
 342 conversion factor between risk (as measured by probability) and a physical variable, zonal wind.

343 An observation of $U(30 \text{ km}) = 38 \text{ m/s}$ implies a 50% chance of vortex breakdown. The variation
 344 in slope also tells us that a wind reduction from 40 m/s to 30 m/s represents a far greater increase
 345 in risk than a reduction from 30 m/s to 20 m/s. Meanwhile, panel (b) shows the committor to be
 346 an increasing function of IHF, since SSW is associated with large wave amplitude and phase lag.
 347 However, IHF seems inferior to zonal wind as a committor proxy, as a small change in IHF from
 348 ~ 0.005 to ~ 0.01 corresponds to a sharp increase in committor from nearly zero to nearly one.
 349 In other words, knowing only IHF doesn't provide much useful information about the threat of
 350 SSW until it is already virtually certain. The dotted envelope is also wider in panel (b) than (a),
 351 indicating that projecting the committor onto IHF removes more information than projecting onto
 352 U . While the underlying noise makes it impossible to divine the outcome with certainty from *any*
 353 observation, the projection error clearly privileges some observables over others for their predictive
 354 power.

355 In panels (c) and (d), the lead time is seen to have the opposite overall trend as the committor: the
 356 weaker the wind, or the greater the heat flux, the closer you are on average to a vortex breakdown.
 357 $\eta^+(\mathbf{x})$ is not defined when wind is strongest, as $\mathbf{x} \in A$ and so $q^+(\mathbf{x}) = 0$. However, an interesting
 358 exception to the trend occurs in the range $10 \text{ km} \leq U \leq 40 \text{ km}$: the expected lead time stays constant
 359 or slightly *decreases* as zonal wind increases, and the projection error remains large. This means
 360 that while the probability of vortex breakdown increases rapidly from 50% to 90%, the time until
 361 vortex breakdown remains highly uncertain. To resolve this seeming paradox, we will have to
 362 visualize the joint variation of q^+ and η^+ .

363 It is of course better to consider multiple observables at once. Figure 3 shows the information
 364 gained beyond observing $U(30 \text{ km})$ by incorporating IHF as a second observable. In the top row we
 365 project π , q^+ , and η^+ onto the two-dimensional subspace, revealing structure hidden from view in
 366 the one-dimensional projections. Panel (a) is a 2-dimensional extension of Figure 1(d), with density

367 peaks visible in the neighborhoods of **a** and **b**. The white space surrounding the gray represents
 368 physically insignificant regions of state space that was not sampled by the long simulation. The
 369 same convention holds for the following two-dimensional figures. The committor is displayed in
 370 panel (b) over the same space. It changes from blue at the top (an SSW is unlikely) to red at
 371 the bottom (an SSW is likely), bearing out the negative association between U and q^+ . However,
 372 there are non-negligible horizontal gradients that show that IHF plays a role, too. Likewise, the
 373 lead time in panel (c) decreases from ~ 90 days near **a** to 0 days near **b**, when the transition is
 374 complete. Here, IHF appears even more critically important for forecasting how the event plays
 375 out, as gradients in η^+ are often completely horizontal.

376 A horizontal dotted line in Figure 3(a-c) marks the 50% risk level $U(30 \text{ km}) = 38 \text{ m/s}$, but the
 377 committor varies along it from low risk at the left to high risk at the right: we show this concretely
 378 by selecting two points θ_0 and θ_1 along the line. According to U alone, i.e., the curve in Figure
 379 2(a), both would have the same committor of 0.5. According to both U and IHF together, i.e., the
 380 two-dimensional heat map in Figure 3(b), they have very different probabilities of $q^+(\theta_0) = 0.37$
 381 and $q^+(\theta_1) = 0.65$: an SSW is nearly twice as likely to occur from starting point θ_1 as θ_0 .

382 While those committor values come from the DGA method to be described in Section 5, we
 383 confirm them empirically by plotting an ensemble of 100 trajectories originating from each of the
 384 two initial conditions in panels (d) and (e) below, coloring A -bound trajectories blue and B -bound
 385 trajectories red. Only 30% of the sampled trajectories through θ_0 exhibit an SSW, next going to
 386 state B , while 63% of the integrations from θ_1 end at B . In both cases, the heatmaps and ensemble
 387 sample means roughly match. The small differences between the projected committor and the
 388 empirical “success” rate of trajectories arises both from errors in the DGA calculation (which we
 389 analyze in section 5) and the finite size of the ensemble.

390 The lead time prediction is improved similarly by incorporating the second observable. According
391 to U alone, Figure 2 predicts a lead time of 40 days for both θ_0 and θ_1 . Considering IHF additionally,
392 the two-dimensional heat map in Figure 3 predicts a lead time of 47 days and 28 days for θ_0 and
393 θ_1 , respectively. Referring to the ensemble from θ_1 in panels (d) and (e), the arrival times of red
394 trajectories to B provide a discrete sampling of the lead time distributions of $\tau_B | \tau_B < \tau_A$. The
395 sample means are 50 and 32 days respectively from θ_0 and θ_1 , again roughly matching with our
396 predictions.

397 These two-dimensional projections still leave out 73 remaining dimensions, which we could
398 incorporate to make the forecasts even better. After accounting for all 75 dimensions, we would
399 obtain the full committor function $q^+ : \mathbb{R}^d \rightarrow \mathbb{R}$. This is still a probability, i.e., an expectation over
400 the unresolved turbulent processes and uncertain initial condition. Low-dimensional committor
401 projections simply treat the projected-out dimensions as random variables sampled according to
402 π . Whether projected to a space of 1 or 75 dimensions, the committor is the function of that space
403 that is closest, in the mean-square sense, to the binary indicator $\mathbb{1}_B(\mathbf{X}(\tau))$; this is the defining
404 characteristic of conditional expectation (Durrett 2013). In the case that the system does hit B next,
405 the lead time is closest in the mean-square sense to τ_B .

406 While high-dimensional systems offer many coordinates to choose from, we argue that the
407 committor and lead time are the most important nonlinear coordinates to monitor for forecasting
408 purposes. We will explore their relationship in the next subsection. Although both encode some
409 version of proximity to SSW, they are independent variables which deserve separate consideration.

410 *d. Relationship between risk and lead time*

411 A forecast is most useful if it comes sufficiently early (to leave some buffer time before impact)
412 and is sufficiently precise to time your response. For example, in June we can say with certainty it

413 will snow next winter in Minnesota. To be useful, we want to know the date of the first snow as
414 early as possible. By relating levels of risk (quantified by q^+) and lead time (quantified by η^+), we
415 can now assess the limits of early prediction. Such a relationship would answer two questions: for
416 an SSW transition, (1) how far in advance will we be aware of it with some prescribed confidence,
417 say 80%? (2) given some prescribed lead time, say 42 days, how aware or ignorant could we be of
418 it?

419 The committor and lead time have an overall negative relationship, but they do not completely
420 determine each other, as the contours in Figure 3(a,b) do not perfectly line up. We treat them as
421 independent variables in Figure 4, which maps zonal wind and IHF as functions of the coordinates
422 q^+ and η^+ in an inversion of Figure 3. The density $\pi(\mathbf{x})$ projected on this space in 4(a) shows again
423 a bimodal structure around **a** and **b**, which occupy opposite corners of this space by construction.
424 Meanwhile, zonal wind and IHF are indicated by the shading in panels (b) and (c). The bridge
425 between **a** and **b** is not a narrow band, but rather includes a curious high-committor, high-lead
426 time branch which seems paradoxical: points at $q^+ = 0.9$ have a greater spread in η^+ than points
427 at $q^+ = 0.5$, contrary to the intuition that closeness to B in probability means closeness in time.
428 The color shading shows that q^+ is strongly associated with $U(30 \text{ km})$, while η^+ is more strongly
429 associated IHF(30 km). In particular the horizontal contours in panel (c) show that the large spread
430 in lead time near B is due almost completely to variation in IHF. In other words, the system can
431 be highly committed to B with a low zonal wind, but if IHF is low, it may take a long time to get
432 there. We can also see this from the lower-left region of Figure 3(a) and (b), where committor is
433 high and lead time is high.

434 There are two complementary explanations for this phenomenon. First, the low- U , low-IHF
435 region of state space corresponds to a temporary restoration phase in a vacillation cycle, which
436 delays the inevitable collapse of zonal wind below the threshold defining B . In fact, the ensemble

437 of pathways starting from θ_0 in Figure 3(c) has one member whose zonal wind repeatedly dips
438 low, but not quite to the level of \mathbf{b} , and partially restores before finally plunging all the way down.
439 These cycles are reminiscent of minor warmings preceding major ones.

440 The second explanation is that many of these partial restoration events are not part of an $A \rightarrow B$
441 transition, but rather a $B \rightarrow B$ transition. In a highly irreversible system such as the Holton-Mass
442 model, these two situations are quite dynamically distinct. To distinguish them using DGA, we
443 would have to account for the *past* as well as the future, calculating backward-in-time forecasts
444 such as the backward committor $q^-(\mathbf{x}) = \mathbb{P}_{\mathbf{x}}\{\mathbf{X}(\tau^-) \in A\}$, where $\tau^- < 0$ is the most-recent hitting
445 time. Backward forecasts will be analyzed thoroughly in a forthcoming paper, but they are beyond
446 the scope of the present one.

447 In summary, q^+ and η^+ are principled metrics to inform preparation for extreme weather. For
448 example, a threatened community might decide in advance to start taking action when an event is
449 very likely, $q^+ \geq 0.8$, and somewhat imminent, $\eta^+ \leq 10$ days, or rather, when an event is somewhat
450 likely, $q^+ \geq 0.5$, and very imminent, $\eta^+ \leq 3$ days. Because of partial restoration events, the
451 committor does not determine the lead time or vice versa, and so a good real-time disaster response
452 strategy should take both of them into account, defining an “alarm threshold” that is not a single
453 number, but some function of both the committor and lead time. This idea is similar in spirit to
454 that of the Torino scale, which assigns a single risk metric to an asteroid or comet impacts based
455 on both probability and severity (Binzel 2000). Of course, after many near-SSW events, a lot of
456 material damage may have already occurred, which may be a reason to define a higher threshold
457 for the definition of B , or even a continuum for different severity levels of SSW. We emphasize
458 that the choice of A , B and alarm thresholds are more of a community and policy decision than a
459 scientific one. The strength of our approach is that it provides a flexible numerical framework to
460 quantify and optimize the consequences of those decisions.

461 **4. Sparse representation of the committor**

462 The committor projections showed give only an impression of its high-dimensional structure.
 463 While Equation (15) says how to optimally represent the committor over a given CV subspace,
 464 optimizing $S[f; \theta]$ over f , it does not say which subspace θ is optimal. If the committor does admit
 465 a sparse representation, we could specifically target observations on these high-impact signals. In
 466 this section we address this much harder problem of optimizing $S[f; \theta]$ over subspaces θ .

467 The set of CV spaces is infinite, as observables θ can be arbitrarily complex nonlinear functions
 468 of the basic state variables \mathbf{x} . Machine learning algorithms such as artificial neural networks
 469 are designed exactly for that purpose: to represent functions nonparametrically from observed
 470 input-output pairs. However, to keep the representation interpretable, we will restrict ourselves
 471 to physics-informed input features based on the Eliassen-Palm (EP) relation, which relates wave
 472 activity, PV fluxes and gradients, and heating source terms in a conservation equation. From Yoden
 473 (1987b), the EP relation for the Holton-Mass model takes the form

$$\begin{aligned} \partial_t \left(\frac{q'^2}{2} \right) + (\partial_y \bar{q}) \rho_s^{-1} \nabla \cdot \mathbf{F} \\ = - \frac{f_0^2}{N^2} \rho_s^{-1} \overline{q' \partial_z (\alpha \rho_s \partial_z \psi')} \end{aligned} \quad (17)$$

$$\text{where } \mathbf{F} = (-\rho_s \overline{u'v'}) \mathbf{j} + (\rho_s \overline{v' \partial_z \psi'}) \mathbf{k}$$

474 The EP flux divergence has two alternative expressions: $\rho_s^{-1} \nabla \cdot \mathbf{F} = \overline{v'q'} = \rho_s^{-1} \frac{R}{Hf_0} \partial_z [\rho_s \overline{v'T'}]$. If
 475 there were no dissipation ($\alpha = 0$) and the background zonal state were time-independent ($\partial_t \bar{q} = 0$),
 476 dividing both sides by $\partial_y \bar{q}$ would express local conservation of wave activity $\mathcal{A} = \rho_s \overline{q'^2} / (2\partial_y \bar{q})$.
 477 Neither of these is exact in the stochastic Holton-Mass model, so we use the quantities in Equation
 478 (17) as diagnostics: enstrophy $\overline{q'^2}$, PV gradient $\partial_y \bar{q}$, PV flux $\overline{v'q'}$, and heat flux $\overline{v'T'}$. Each field is
 479 a function of (y, z) and takes on very different profiles for the states \mathbf{a} and \mathbf{b} , as found by Yoden
 480 (1987b). A transition from A to B , where the vortex weakens dramatically, must entail a reduction

481 in $\partial_y \bar{q}$ and a burst in positive $\overline{v'T'}$ (negative $\overline{v'q'}$) as a Rossby wave propagates from the tropopause
 482 vertically up through the stratosphere and breaks. This is the general physical narrative of a sudden
 483 warming event, and these same fields might be expected to be useful observables to track for
 484 qualitative understanding and prediction. For visualization, we have found $U(30\text{ km})$ and $\text{IHF}(30$
 485 $\text{ km}) = \int_{0\text{ km}}^{30\text{ km}} e^{-z/H} \overline{v'T'} dz$ to be particularly helpful. However, this doesn't necessarily imply they
 486 are optimal predictors of q^+ , and regression is a more principled way to find them.

487 We start by projecting the committor onto each observable at each altitude separately, in hopes
 488 of finding particularly salient altitude levels that clarify the role of vertical interactions. The first
 489 five rows of Figure 5 display, for five fields (U , $|\Psi|$, $\overline{q'^2}$, $\partial_y \bar{q}$, and $\overline{v'q'}$) and for a range of altitude
 490 levels, the mean and standard deviation of the committor projected onto that field at that altitude.
 491 Each altitude has a different range of the CV; for example, because U has a Dirichlet condition
 492 at the bottom and a Neumann condition at the top, the lower levels have a much smaller range of
 493 variability than the high levels. We also plot the integrated variance, or L^2 projection error, at each
 494 level in the right-hand column. A low projected committor variance over U at altitude z_0 means
 495 that the committor is mostly determined by the single observable $U(z_0)$, while a high projected
 496 variance indicates significant dependence of q^+ on variables other than $U(z_0)$. In order to compare
 497 different altitudes and fields as directly as possible, the L^2 projection error at each altitude is an
 498 average over discrete bins of the observable.

499 In selecting good CV's, we generally look for a simple, hopefully monotonic, and sensitive
 500 relationship with the committor. Of all the candidate fields, U and $\partial_y \bar{q}$ stand out the most in
 501 this respect, being clearly negatively correlated with the forward committor at all altitudes. The
 502 associated projection error tends to be greatest in the region $q^+ \approx 0.5$, as observed before, but
 503 interestingly there is a small altitude band around 15 – 25 km where its magnitude is minimized.
 504 This suggests an optimal altitude for monitoring the committor through zonal wind, giving the

505 most reliable estimate possible for a single state variable. In contrast, the projection of q^+ onto
 506 $|\Psi|$, displays a large variance across all altitudes. The eddy enstrophy and potential vorticity
 507 flux are also rather unhelpful as early warning signs, despite their central role in SSW evolution.
 508 For example, the large, positive spikes in heat flux across all altitudes generally occur after the
 509 committor ≈ 0.5 threshold has already been crossed. Furthermore, the relationship of $\overline{v'q'}$ with the
 510 committor is not smooth. The $q^+ < 0.5$ region at each altitude is a thin band near zero.

511 The exhaustive CV search in Figure 5 is visually compelling in favor of some fields and some
 512 altitudes over others, but it is not satisfactory as a rigorous comparison. Differences between units
 513 and ranges make it difficult to objectively compare the L^2 projection error. Furthermore, restricting
 514 to one variable at a time is limiting. Accordingly, we also perform a more automated approach
 515 to identify salient variables in the form of a generalized linear model for the forward committor,
 516 using sparsity-promoting LASSO regression (“Least Absolute Shrinkage and Selection Operator”)
 517 due to Tibshirani (1996), as implemented in the `scikit-learn` Python package (Pedregosa et al.
 518 2011). As input features, we use all state variables $\text{Re}\{\Psi\}, \text{Im}\{\Psi\}, U$, the integrated heat flux
 519 $\int_0^z e^{-z/H} \overline{v'T'} dz$, the eddy PV flux $\overline{v'q'}$, and the background PV gradient $\partial_y \bar{q}$, at all altitudes z
 520 simultaneously. The advantage of a sparsity-promoting regression is that it isolates a small number
 521 of observables that can accurately approximate the committor in linear combination. Considering
 522 that regions close to A and B have low committor uncertainty, we regress only on data points with
 523 $q^+ \in (0.2, 0.8)$, and of those only a subset weighted by $\pi(\mathbf{x})q^+(\mathbf{x})(1 - q^+(\mathbf{x}))$ to further emphasize
 524 the transition region $q^+ \approx 0.5$. To constrain committor predictions to the range $(0, 1)$, we regress
 525 on the committor after an inverse-sigmoid transformation, $\ln(q^+/(1 - q^+))$. First we do this at each
 526 altitude separately, and in Figure 6 (a) we plot the coefficients of each component as a function of
 527 altitude. The bottom row of Figure 5 also displays the committor projected on the height-dependent
 528 LASSO predictor.

529 The height-dependent regression in 6(a) shows each component is salient for some altitude range.
530 In general, U and $\text{Im}\{\Psi\}$ dominate as causal variables at low altitudes, while $\text{Re}\{\Psi\}$ dominates
531 at high altitudes. The overall prediction quality, as measured by R^2 and plotted in Figure 6 (b), is
532 greatest around 21.5 km, consistent with our qualitative observations of Figure 5. Note that not all
533 single-altitude slices are sufficient for approximating the committor, even with LASSO regression;
534 in the altitude band 50 – 60 km, the LASSO predictor is not monotonic and has a large projected
535 variance, as seen in the bottom row of Figure 5. The specific altitude can matter a great deal. But
536 by using all altitudes at once, the committor approximation may be improved further. We thus
537 repeat the LASSO with all altitudes simultaneously and find the sparse coefficient structure shown
538 in 6 (c), with a few variables contributing the most, namely the state variables Ψ and U in the
539 altitude range 15-22 km. The nonlinear CVs failed to make any nonzero contribution to LASSO,
540 and this remained stubbornly true for other nonlinear combinations not shown, such as $\overline{v'T'}$. With
541 multiple lines of evidence indicating 21.5 km as an altitude with high predictive value for the
542 forward committor, we can make a strong recommendation for targeting observations here. This
543 conclusion applies only to the Holton-Mass model under these parameters, but the methodology
544 explained above can be applied similarly to models of arbitrary complexity.

545 We have presented the committor and lead time as “ideal” forecasts, especially the committor,
546 which we have devoted considerable effort to approximating in this section. We want to emphasize
547 that q^+ and η^+ are not competitors to ensemble forecasting; rather, they are two of its most important
548 end results. So far, we have simply advocated including q^+ and η^+ as quantities of interest. Going
549 forward, however, we do propose an alternative to ensemble forecasting aimed specifically at the
550 committor, lead time, and a wider class of forecasting functions, as they are important enough
551 in their own right to warrant dedicated computation methods. Our approach uses only short
552 simulations, making it highly parallelizable, and shifts the numerical burden from online to offline.

553 Figures 2-6 were all generated using the short-simulation algorithm. While the method is not
 554 yet optimized and in some cases not competitive with ensemble forecasting, we anticipate such
 555 methods will be increasingly favorable with modern trends in computing.

556 5. The computational method

557 In this section we describe the methodology, which involves some technical results from stochastic
 558 processes and measure theory. After describing the theoretical motivation and the numerical
 559 pipeline in turn, we demonstrate the method's accuracy and discuss its efficiency compared to
 560 straightforward ensemble forecasting.

561 a. Feynman-Kac formulae

562 The forecast functions described above—committors and passage times—can all be derived from
 563 general conditional expectations of the form

$$F(\mathbf{x}; \lambda) = \mathbb{E}_{\mathbf{x}} \left[G(\mathbf{X}(\tau)) \exp \left(\lambda \int_0^\tau \Gamma(\mathbf{X}(s)) ds \right) \right] \quad (18)$$

564 where again the subscript \mathbf{x} denotes conditioning on $\mathbf{X}(0) = \mathbf{x}$; G, Γ are arbitrary known functions
 565 over \mathbb{R}^d ; and τ is a stopping time, specifically a first-exit time like Equation (10) but possibly
 566 with D replaced by another set. λ is a variable parameter that turns F into a moment-generating
 567 function. To see that the forward committor takes on this form, set $G(\mathbf{x}) = \mathbb{1}_B(\mathbf{x})$, $\lambda = 0$ (Γ can be
 568 anything), and $\tau = \tau_{A \cup B}$. Then $F(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[\mathbb{1}_B(\mathbf{X}(\tau))] = \mathbb{P}_{\mathbf{x}}\{\mathbf{X}(\tau_{D^c}) \in B\} = q^+(\mathbf{x})$. For the η^+ , set
 569 $\tau = \tau_B$, $G = \mathbb{1}_B$, and $\Gamma = 1$. Then

$$F(\mathbf{x}; \lambda) = \mathbb{E}_{\mathbf{x}}[\mathbb{1}_B(\mathbf{X}(\tau)) \exp(\lambda \tau)] \quad (19)$$

$$\frac{1}{q^+(\mathbf{x})} \frac{\partial}{\partial \lambda} F(\mathbf{x}; 0) = \frac{\mathbb{E}_{\mathbf{x}}[\tau \mathbb{1}_B(\mathbf{X}(\tau))]}{\mathbb{E}_{\mathbf{x}}[\mathbb{1}_B(\mathbf{X}(\tau))]} \quad (20)$$

$$= \eta^+(\mathbf{x}). \quad (21)$$

570 So we must also be able to differentiate F with respect to λ .

571 More generally, the function G is chosen by the user to quantify risk at the terminal time τ ; in
 572 the case of the forward committor, that risk is binary, with an SSW representing a positive risk
 573 and a radiative vortex no risk at all. The function Γ is chosen to quantify the risk accumulated up
 574 until time τ , which might be simply an event's duration, but other integrated risks may be of more
 575 interest for the application. For example, one could express the total poleward heat flux by setting
 576 $\Gamma = \overline{v'T'}$, or the momentum lost by the vortex by setting $\Gamma(\mathbf{x}) = U(\mathbf{a}) - U(\mathbf{x})$. Extending the trick
 577 in (20), one can compute not only means but higher moments of such integrals by expressing the
 578 risk with Γ . Repeated differentiation $F(\mathbf{x}; \lambda)$ gives

$$\partial_\lambda^k F(\mathbf{x}; 0) = \mathbb{E}_{\mathbf{x}} \left[G(\mathbf{X}(\tau)) \left(\int_0^\tau \Gamma(\mathbf{X}(s)) ds \right)^k \right] \quad (22)$$

579 We choose to focus on expectations of the form (18) in order to take advantage of the Feynman-
 580 Kac formula, which represents $F(\mathbf{x}; \lambda)$ as the solution to a PDE boundary value problem over
 581 state space. As PDEs involve local operators, this form is more amenable to solution with short
 582 trajectories which don't stray far from their source. The boundary value problem associated with
 583 (18) is

$$\begin{cases} (\mathcal{L} + \lambda\Gamma)F(\mathbf{x}; \lambda) = 0 & \mathbf{x} \in D \\ F(\mathbf{x}; \lambda) = G(\mathbf{x}) & \mathbf{x} \in D^c \end{cases} \quad (23)$$

584 The domain D here is some combination of A^c and B^c . The operator \mathcal{L} is known as the *infinitesimal*
 585 *generator* of the stochastic process, which acts on functions by pushing expectations forward in
 586 time along trajectories:

$$\mathcal{L}f(\mathbf{x}) := \lim_{\Delta t \rightarrow 0} \frac{\mathbb{E}_{\mathbf{x}}[f(\mathbf{X}(\Delta t))] - f(\mathbf{x})}{\Delta t} \quad (24)$$

587 In a diffusion process like the stochastic Holton-Mass model, \mathcal{L} is an advection-diffusion partial
 588 differential operator which is analogous to a material derivative in fluid mechanics. The generator

589 encapsulates the properties of the stochastic process. In addition to solving boundary value
 590 problems (18), its adjoint \mathcal{L}^* provides the Fokker-Planck equation for the stationary density $\pi(\mathbf{x})$:

$$\mathcal{L}^* \pi(\mathbf{x}) = 0 \quad (25)$$

591 We can also write equations for moments of F , as in (22), by differentiating (23) repeatedly and
 592 setting $\lambda = 0$:

$$\mathcal{L}[\partial_\lambda^k F](\mathbf{x}; 0) = -k\Gamma\partial_\lambda^{k-1}F \quad (26)$$

593 This is an application of the Kac Moment Method (Fitzsimmons and Pitman 1999). Note that we
 594 never actually have to solve (23) with nonzero λ . Instead we implement the recursion above. Note
 595 that the base case, $k = 0$, with $G = \mathbb{1}_B$ gives $F^+ = q^+$, no matter what the risk function Γ . In this
 596 paper we compute only up to the first moment, $k = 1$. Further background regarding stochastic
 597 processes and Feynman-Kac formulae can be found in Karatzas and Shreve (1998); Oksendal
 598 (2003); E et al. (2019).

599 *b. Dynamical Galerkin Approximation*

600 To solve the boundary value problem (23) with $\lambda = 0$, we start by following the standard finite
 601 element recipe, converting to a variational form and projecting onto a finite basis. First, we
 602 homogenize boundary conditions by writing $F(\mathbf{x}) = \hat{F}(\mathbf{x}) + f(\mathbf{x})$, where \hat{F} is a guess function that
 603 obeys the boundary condition $\hat{F}|_{D^c} = G$, and $f|_{D^c} = 0$. Next, we integrate the equation against any
 604 test function ϕ , weighting the integrand by a density μ (which is arbitrary for now, but will be
 605 specified later):

$$\int_{\mathbb{R}^d} \phi(\mathbf{x}) \mathcal{L}f(\mathbf{x}) \mu(\mathbf{x}) d\mathbf{x} = \int \phi(\mathbf{x}) (G - \mathcal{L}\hat{F})(\mathbf{x}) \mu(\mathbf{x}) d\mathbf{x}$$

$$\langle \phi, \mathcal{L}f \rangle_\mu = \langle \phi, G - \mathcal{L}\hat{F} \rangle_\mu \quad (27)$$

606 The test function ϕ should live in the same space as f , i.e., with homogeneous boundary conditions
607 $\phi(\mathbf{x}) = 0$ for $\mathbf{x} \in A \cup B$. We refer to the inner products in (27) as being “with respect to” the
608 measure (with density) μ . We approximate f by expanding in a finite basis $f(\mathbf{x}) = \sum_{j=1}^M \xi_j \phi_j(\mathbf{x})$
609 with unknown coefficients ξ_j , and enforce that (27) hold for each ϕ_i . This reduces the problem to
610 a system of linear equations,

$$\sum_{j=1}^M \langle \phi_i, \mathcal{L}\phi_j \rangle_{\mu} \xi_j = \langle \phi_i, G - \mathcal{L}\hat{F} \rangle_{\mu} \quad i = 1, \dots, M \quad (28)$$

611 which can be solved with standard numerical linear algebra packages.

612 This procedure consists of three crucial subroutines. First, we must construct a set of basis
613 functions ϕ_j . Second, we have to evaluate the generator’s action on them, $\mathcal{L}\phi_j$. Third, we have
614 to compute inner products. With standard PDE methods, the basis size would grow exponentially
615 with dimension, quickly rendering the first and third steps intractable. Successful approaches
616 will involve a representation of the solution, F , suitable for the high dimensional setting, i.e.
617 representations of the type commonly employed for machine learning tasks. DGA is one such
618 method, whose special twist is to construct a “data-informed” basis of reasonable size, evaluate
619 the generator by implementing Equation (24) with the same data set, and finally evaluate the inner
620 products (27) with a Monte Carlo integral. The data consist of short trajectories launched from
621 all over state space, which the system of linear equations stitches together into a global function
622 estimate. We sketch the procedure here, but for the implementation details we refer to the appendix
623 and to Thiede et al. (2019) and Strahan et al. (2021), where DGA has already been developed for
624 molecular dynamics.

625 **Step 1:** Generate the data, in the format of N initial conditions $\{\mathbf{X}_n : 1 \leq n \leq N\}$. Evolve each
626 initial condition forward for a “lag time” Δt to obtain a set of short trajectories $\{\mathbf{X}_n(t) : 0 \leq t \leq$
627 $\Delta t, n = 1, \dots, N\} \subset \mathbb{R}^d$. Here and going forward, \mathbf{X}_n will mean $\mathbf{X}_n(0)$. The choice of starting

628 points is flexible, but crucial for the efficiency and accuracy of DGA. Because our goal here is to
 629 demonstrate interpretable results, we prioritize simplicity and accuracy over efficiency, and defer
 630 optimization to later work. We simply draw initial conditions at random from the long control
 631 simulation of 5×10^5 days, and then generate new short trajectories from those points. We do not
 632 sample the points with equal probability, but instead re-weight to get a uniform distribution over
 633 the space $(U(30\text{km}), |\Psi|(30\text{km}))$, within the bounds realized by the control simulation, which
 634 are approximately $-30\text{ m/s} \leq U(30\text{km}) \leq 70\text{ m/s}$ and $0\text{ m}^2/\text{s} \leq |\Psi|(30\text{km}) \leq 2 \times 10^7\text{ m}^2/\text{s}$. This
 635 sampling procedure, and any other version, implicitly defines a *sampling measure* μ on state
 636 space, where $\mu(\mathbf{x}) d\mathbf{x}$ is the expected fraction of starting points in the neighborhood $d\mathbf{x}$ about \mathbf{x} .
 637 Sampling points with equal weight from the control run would induce $\mu = \pi$, a very inefficient
 638 choice because probability concentrates around the metastable states **a** and **b**. The re-weighting
 639 procedure ensures data coverage of intermediate-wind regions between *A* and *B*, as well as the
 640 large bursts of wave amplitude that characterize the transition pathways. Our main results use
 641 $N = 5 \times 10^5$ short trajectories with a lag time of $\Delta t = 20$ days, sampled at a frequency of twice per
 642 day. This data set is more than needed to get a reasonable committor estimate, but we have sampled
 643 generously in order to visualize the functions in high detail. The final section will show the method
 644 is robust, capable of reasonably approximating the committor even with an order-of-magnitude
 645 reduction in data.

646 **Step 2:** Define the basis. The Galerkin method works for any class of basis functions that becomes
 647 increasingly expressive as the library grows and becomes capable of estimating any function of
 648 interest. However, with a finite truncation, choosing basis functions is a crucial ingredient of DGA,
 649 greatly impacting the efficiency and accuracy of the results. In our current study, we restrict to the
 650 simplest kind of basis, which consists of indicator functions $\phi_i(x) = \mathbb{1}_{S_i}(x)$, where $\{S_1, \dots, S_M\}$ is a

651 disjoint partition of state space. In practice we will construct these sets by clustering the initial data
 652 points as described in more detail in Appendix A. This is a common practice in the computational
 653 statistical mechanics community for building a Markov State Model (MSM) (Chodera et al. 2006;
 654 Frank and Fischer 2008; Pande et al. 2010; Bowman et al. 2013; Chodera and Noé 2014). MSMs
 655 are a dimensionality reduction technique that has also been used in conjunction with analysis of
 656 metastable transitions, primarily in protein folding dynamics (Noé et al. 2009). MSMs have also
 657 been used recently to study garbage patch dynamics in the ocean (Miron et al. 2021) as well as
 658 complex social dynamics (Helfmann et al. 2021). In Maiocchi et al. (2020), the authors take
 659 an interesting approach to MSMs by clustering points based on proximity to unstable periodic
 660 orbits, a potentially useful paradigm for general chaotic weather phenomena (Lucarini and Gritsun
 661 2020). DGA can be viewed as an extension of MSMs, though, rather than producing any reduced
 662 complexity model, the explicit goal in DGA is estimating specific functions as in Equation (18).

663 **Step 3:** Apply the generator. The forward difference formula

$$\widehat{\mathcal{L}}\phi(\mathbf{X}_n) = \frac{\phi(\mathbf{X}_n(\Delta t)) - \phi(\mathbf{X}_n)}{\Delta t} \quad (29)$$

664 suggested by the definition of the generator (24), results in a systematic bias when Δt is finite. On
 665 the other hand, small values of Δt lead to large variances in our Monte Carlo estimates of the inner
 666 products in (28). To resolve these issues we use an integrated form of the Feynman–Kac equations
 667 that involves stopping trajectories when they enter A or B . Details are provided in Appendix A.

668 **Step 4:** Compute the inner products. The inner products in Equation (28) are integrals over high-
 669 dimensional state space that are intractable with standard quadrature, but can be approximated
 670 using Monte Carlo integration. If \mathbf{X} is an \mathbb{R}^d -valued random variable distributed according to μ ,
 671 and we have access to random samples $\{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ (which we do), the law of large numbers

672 gives, for any function g with finite expectation,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N g(\mathbf{X}_n) = \int_{\mathbb{R}^d} g(\mathbf{x}) \mu(\mathbf{x}) d\mathbf{x} \quad (30)$$

673 Setting $g(\mathbf{x}) = \phi_i(\mathbf{x}) \mathcal{L} \phi_j(\mathbf{x})$, the sample average on the left-hand side of (30) therefore provides
 674 an estimator of $\langle \phi_i, \mathcal{L} \phi_j \rangle$. Of course, our approximation uses finite N and nonzero Δt . A similar
 675 sample average approximation can be used to estimate the inner product on the right-hand side of
 676 (28).

677 These same steps apply to both q^+ and $\mathbb{E}[\tau_B]$, as well as the recursion in (26) for η^+ . For the
 678 Fokker-Planck equation (25), one extra step is needed to convert an equation with \mathcal{L}^* into an
 679 equation with \mathcal{L} . Our procedure for estimating π is described in Appendix A.

680 **Step 5:** Solve the equation (28). With a reasonable basis size $M \lesssim 1000$, an LU solver such as in
 681 LAPACK via Numpy can handle Equation (28). In the case of the homogeneous system for $w(\mathbf{x})$,
 682 a QR decomposition can identify the null vector.

683 *c. DGA fidelity and sensitivity analysis*

684 To illustrate the effect of parameter choices on performance, we present here a simple sensitivity
 685 analysis. Figure 7 verifies the numerical accuracy and convergence of DGA by plotting the
 686 committor as a function of $U(30 \text{ km})$, estimated both from the control simulation and with DGA,
 687 for various DGA parameters. The red curves $q_{\text{DGA}}^+(U(30 \text{ km}))$ are calculated by projecting the
 688 committor as in Figure 2(a), while the black curve $q_{\text{EMP}}^+(U(30 \text{ km}))$ is an empirical committor
 689 estimate equal to the fraction of control simulation points seen at a particular value of $U(30 \text{ km})$
 690 that next hit B .

691 In panels (a), (b), and (d), the lag time Δt increases from 5 to 10 to 20 days while the number
 692 of short trajectories stays fixed at $N = 5 \times 10^5$. Panel (c) has a long lag of 20 days, but a small

693 data set of $N = 5 \times 10^4$, allowing us to see the tradeoff between N and Δt . The basis size M is
 694 chosen heuristically as large as possible within reason for the clustering algorithm (see Appendix
 695 A). While DGA tends to systematically overestimate q^+ relative to q_{EMP}^+ in the mid-range of
 696 U , it seems to approach the empirical estimate as the data size and lag time increase. Each
 697 plot also displays the root-mean-square deviation between the two estimators over this subspace,
 698 $\varepsilon = \sqrt{\langle (q_{\text{DGA}}^+ - q_{\text{EMP}}^+)^2 \rangle_\pi}$. Within this regime, it seems that increasing the lag time has a greater
 699 impact on the deviation than increasing the number of data points. Panels (b) and (c) have
 700 approximately the same deviation ε , but (c) uses only one fifth the data, measured by total
 701 simulation time. On the other hand, more short trajectories can be parallelized more readily than
 702 fewer long trajectories, and the optimal choice will depend on computing resources.

703 It is natural to ask whether our short trajectory based approach is more efficient than a direct
 704 shooting approach in which many independent “long” trajectories are launched from a single initial
 705 condition \mathbf{x} and the committor probability $q^+(\mathbf{x})$ (or another forecast) is estimated directly. For a
 706 single value of \mathbf{x} for which $q^+(\mathbf{x})$ is not very small (so that a non-negligible fraction of trajectories
 707 reach B before A) and for which the lead time $\eta^+(\mathbf{x})$ is not too large (so that trajectories reaching B
 708 do so without requiring long integration times), direct shooting will undoubtedly be more efficient.
 709 However, a key feature of our approach is that it simultaneously estimates forecasts at all values
 710 of \mathbf{x} , allowing the subsequent analysis of those functions that has been the focus of much of this
 711 article. Building accurate estimators in all of state space by direct shooting would be extremely
 712 costly even for the reduced complexity model studied here.

713 6. Conclusion

714 Forecasting rare events is, by the very nature of rare events, an extremely difficult computational
 715 task, and one of science’s most pressing challenges. We have described a computational framework,

716 a dynamical Galerkin approximation to the Feynman-Kac equations, that combines the minimalistic
717 philosophy of dimensionality reduction with the fidelity of high-resolution models. We identify
718 a set of reduced coordinates, the committor probability and expected lead time, that provide the
719 essential information that large ensemble forecasts hope to compute. DGA uses relatively short
720 simulations of the full model to estimate these quantities of interest, allowing for prediction on
721 much longer timescales than that of the simulation. In its focus on directly estimating statistics
722 of interest, DGA differs from previous reduced-order modeling methods that attempt to capture
723 general qualities of the system, including both physics-based models (Lorenz 1963; Charney and
724 DeVore 1979; Legras and Ghil 1985; Crommelin 2003; Timmermann et al. 2003; Ruzmaikin et al.
725 2003) and more recent data-driven models making use of machine learning (Giannakis and Majda
726 2012; Giannakis et al. 2018; Berry et al. 2015; Sabeerali et al. 2017; Majda and Qi 2018; Wan et al.
727 2018; Bolton and Zanna 2019; Chattopadhyay et al. 2020; Chen and Majda 2020; Kashinath et al.
728 2021; Chattopadhyay et al. 2021).

729 We have shown numerical results in the context of a stochastically forced Holton-Mass model with
730 75 degrees of freedom, which points to the method's promise for forecasting. By systematically
731 evaluating many model variables for their utility in predicting the fate of the vortex, we have
732 identified some salient physical descriptions of early warning signs. We have furthermore examined
733 the relationship between probability and lead time for a given rare event, a powerful pairing for
734 assessing predictability and preparing for extreme weather. Our results suggest that the slow
735 evolution of vortex preconditioning is an important source of predictability. In particular, the zonal
736 wind and streamfunction in the range of 10-20 km above the tropopause seems to be optimal among
737 a large class of dynamically motivated observables.

738 Beyond the problem of real-time weather forecasting, it is also important to assess the climatology,
739 i.e., long-term frequency, intensity, and other characteristics of rare events. For this goal as well,

740 our methodology offers advantages over large ensemble simulations, which are currently the most
741 detailed source of data (e.g., Schaller et al. 2018). The committor and lead time are ingredients
742 in a larger framework called Transition Path Theory (TPT) for describing rare transition events *at*
743 *steady state*, meaning average properties over long timescales. TPT describes not only the future
744 evolution from an initial condition ($\mathbf{x} \rightarrow B$), but the ensemble of full vortex breakdown events
745 ($A \rightarrow B$), and how they differ from restoration events ($B \rightarrow A$). In principle, interrogating the
746 ensemble of transition paths requires direct simulation of the system long enough to observe many
747 transition events. However, using TPT, quantities computable by our framework can be combined
748 to yield key statistics describing the ensemble of transition paths (Metzner et al. 2006, 2009;
749 Vanden-Eijnden and E 2010; E. and Vanden-Eijnden 2006; Finkel et al. 2020). In a following
750 paper we will apply the same short-trajectory forecasting approach together with TPT to compute
751 transition path statistics such as return times and extract insight about physical mechanisms of the
752 transition process.

753 Scaling our approach up to state-of-the-art weather and climate models will require significant
754 further development. In particular, a completely new procedure for generating trajectory initial
755 conditions will need to be introduced. Generation of a trajectory long enough to thoroughly
756 sample transitions will not be practical for more complicated models. One promising alternative
757 is launching many trajectories in parallel and selectively replicating those that explore new regions
758 of state space, especially transition regions. Such an approach could build on exciting progress
759 over the last decade in targeted rare event simulation schemes (Hoffman et al. 2006; Weare 2009;
760 Bouchet et al. 2011, 2014; Vanden-Eijnden and Weare 2013; Chen et al. 2014; Yasuda et al. 2017;
761 Farazmand and Sapsis 2017; Dematteis et al. 2018; Mohamad and Sapsis 2018; Dematteis et al.
762 2019; Webber et al. 2019; Bouchet et al. 2019a,b; Plotkin et al. 2019; Simonnet et al. 2020; Ragone
763 and Bouchet 2020; Sapsis 2021). A potential challenge here is that GCMs may not be set up for

764 short simulations that start and stop frequently. For this reason, it may be sensible to use longer lag
765 times and a sliding window to define short trajectories. Defining the source of stochasticity is also
766 an important step that varies between models. Explicitly stochastic parameterization (e.g., Berner
767 et al. 2009; Porta Mana and Zanna 2014) will automatically lead to a spread in the short-trajectory
768 ensemble, but in deterministic models, uncertainty will arise from perturbing the initial conditions.
769 This may require special care depending on the model.

770 Another area of algorithmic improvement is selecting a basis expansion of the forecast functions.
771 In upcoming work we will explore more flexible representations using kernel methods and neural
772 networks. The solution of high-dimensional PDEs is an active research area that is making
773 innovative use of machine learning, particularly in the fields of computational chemistry, quantum
774 mechanics, and fluid dynamics (e.g., Carleo and Troyer 2017; Han et al. 2018; Khoo et al. 2018; Li
775 et al. 2020; Mardt et al. 2018; Li et al. 2019; Raissi et al. 2019; Lorpaiboon et al. 2020). Similar
776 approaches may hold great potential for understanding predictability in atmospheric science.

777 *Acknowledgments.* J.F. is supported by the U.S. Department of Energy, Office of Science, Office of
778 Advanced Scientific Computing Research, Department of Energy Computational Science Graduate
779 Fellowship under Award Number DE-SC0019323. R.J.W. is supported by New York University’s
780 Dean’s Dissertation Fellowship and by the Research Training Group in Modeling and Simulation
781 funded by the National Science Foundation via grant RTG/DMS-1646339. E.P.G. acknowledges
782 support from the U. S. National Science Foundation through grant AGS-1852727. This work was
783 partially supported by the NASA Astrobiology Program, grant No. 80NSSC18K0829 and benefited
784 from participation in the NASA Nexus for Exoplanet Systems Science research coordination
785 network. J.W. acknowledges support from the Advanced Scientific Computing Research Program
786 within the DOE Office of Science through award DE-SC0020427. The computations in the paper

787 were done on the high-performance computing clusters at New York University and the Research
788 Computing Center at the University of Chicago.

789 We extend special thanks to Thomas Birner and two anonymous reviewers from *Monthly Weather*
790 *Review*, who provided invaluable feedback on both the technical and high-level aspects of the
791 manuscript. Their insight has helped us to sharpen and clarify the message. We thank John
792 Strahan, Aaron Dinner, and Chatipat Lorpaiboon for helpful methodological advice. Mary Silber,
793 Noboru Nakamura, and Richard Kleeman offered invaluable scientific insight. J.F. benefitted from
794 many helpful discussions with Anya Katsevich.

795 *Data availability statement.* The code for simulating the model, performing DGA, and producing
796 plots is publicly available in the SHORT Github repository, “Solving for Harbingers Of Rare
797 Transitions”, at <https://github.com/justinfoocus12/SHORT>. J.F. is happy to provide further guidance
798 upon request.

799 APPENDIX A

800 **Feynman-Kac formula and DGA**

801 In this section we spell out the DGA procedure in more detail than the main text, explaining the
802 variants that get us to the more intricate conditional expectations. The theoretical background
803 can be found in, e.g., Karatzas and Shreve (1998); Oksendal (2003); E et al. (2019). Let $\mathbf{X}(t)$
804 be a time-homogeneous stochastic process with continuous sample paths in \mathbb{R}^d . Associated to
805 this process is the infinitesimal generator, \mathcal{L} , which acts on observable functions by evolving their
806 expectation forward in time:

$$\mathcal{L}f(\mathbf{x}) = \lim_{\Delta t \rightarrow 0} \frac{\mathbb{E}_{\mathbf{x}}[f(\mathbf{X}(\Delta t))] - f(\mathbf{x})}{\Delta t} \quad (\text{A1})$$

807 where $\mathbb{E}_{\mathbf{x}}[\cdot] := \mathbb{E}[\cdot | \mathbf{X}(0) = \mathbf{x}]$. It can be shown that under the above assumptions on \mathbf{X} , the Itô chain
 808 rule gives

$$df(\mathbf{X}(t)) = \mathcal{L}f(\mathbf{X}(t)) dt + d\mathbf{M}(t) \quad (\text{A2})$$

809 where $\mathbf{M}(t)$ is a martingale. More concretely, in this paper, $\mathbf{X}(t)$ is an Itô diffusion obeying the
 810 stochastic differential equation

$$\begin{aligned} \mathbf{X}(t) = \mathbf{X}(0) &+ \int_0^t b(\mathbf{X}(s)) ds \\ &+ \int_0^t \sigma(\mathbf{X}(s)) d\mathbf{W}(s) \end{aligned} \quad (\text{A3})$$

811 with infinitesimal generator and martingale terms

$$\begin{aligned} \mathcal{L}f(\mathbf{x}) = &\sum_{i=1}^d b_i(\mathbf{x}) \frac{\partial f(\mathbf{x})}{\partial x_i} \\ &+ \sum_{i=1}^d \sum_{j=1}^d \frac{1}{2} [\sigma(\mathbf{x}) \sigma(\mathbf{x})^\top]_{ij} \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} \end{aligned} \quad (\text{A4})$$

$$d\mathbf{M}(t) = \sum_{i=1}^d \frac{\partial f(\mathbf{x})}{\partial x_i} \sigma_{ij}(\mathbf{x}) d\mathbf{W}_j(t) \quad (\text{A5})$$

812 The key forecasting quantities in this paper are of the form (18) and can be solved with (23), a
 813 linear equation involving the generator. We now lay out a brief derivation of the Feynman-Kac
 814 formula and our numerical discretization, roughly following E et al. (2019).

815 *a. Feynman-Kac formula*

816 Let D be a domain in \mathbb{R}^d (for example, $(A \cup B)^c$) and $\tau_{D^c} = \min\{t \geq 0 : \mathbf{X}(t) \notin D\}$ be the first
 817 exit time from this domain starting at time zero. This is a random variable which depends on the
 818 starting condition $\mathbf{x} \in D$. Let $G : \partial D \rightarrow \mathbb{R}$ be a boundary condition, $\Gamma : D \rightarrow \mathbb{R}$ a source term, and
 819 $\Gamma : D \rightarrow \mathbb{R}$ a term to represent accumulated risk. We seek a PDE for the conditional expectation

820 from (18):

$$F(\mathbf{x}) = \mathbb{E}_{\mathbf{x}} \left[G(\mathbf{X}(\tau)) \exp \left(\lambda \int_0^\tau \Gamma(\mathbf{X}(s)) ds \right) \right] \quad (\text{A6})$$

821 where $\mathbb{E}_{\mathbf{x}}[\cdot] = \mathbb{E}[\cdot | \mathbf{X}(0) = \mathbf{x}]$. To derive the PDE (23), consider the following stochastic process:

$$Z(t) = F(\mathbf{X}(t))Y(t) \quad (\text{A7})$$

822 where $Y(t) := \exp \left(\lambda \int_0^t \Gamma(\mathbf{X}(s)) ds \right)$. Itô's lemma gives us that $dY(t) = \lambda \Gamma(\mathbf{X}(t))Y(t) dt$. Hence,

823 applying the product rule to $Z(t)$,

$$dZ(t) = dF(\mathbf{X}(t))Y(t) + F(\mathbf{X}(t)) dY(t) \quad (\text{A8})$$

$$= \mathcal{L}F(\mathbf{X}(t))Y(t) dt + d\mathbf{M}(t)Y(t) \quad (\text{A9})$$

$$+ \lambda F(\mathbf{X}(t))\Gamma(\mathbf{X}(t))Y(t) dt$$

$$= [\mathcal{L}F + \lambda \Gamma F](\mathbf{X}(t))Y(t) dt + Y(t)d\mathbf{M}(t) \quad (\text{A10})$$

824 where in (A8) we have left out the quadratic cross-variation of $F(\mathbf{X}(t))$ and $Y(t)$ because Y has

825 finite variation. If the bracketed term $(\mathcal{L} + \lambda \Gamma(\mathbf{x}))F(\mathbf{x}) = 0$ for all \mathbf{x} , then $Z(t)$ is a martingale and

826 it follows that

$$Z(0) = \mathbb{E}_{\mathbf{x}}[Z(t)] \quad (\text{A11})$$

$$F(\mathbf{x}) = \mathbb{E}_{\mathbf{x}} \left[F(\mathbf{X}(t)) \exp \left(\lambda \int_0^t \Gamma(\mathbf{X}(s)) ds \right) \right] \quad (\text{A12})$$

827 Finally, the formula still holds if we substitute a stopping time for t . By choosing τ , the first exit

828 time from D , the $F(\mathbf{X}(t))$ inside the brackets becomes its boundary value $G(\mathbf{X}(\tau))$. Thus $F(\mathbf{x})$ as

829 defined in (A6) also solves the PDE boundary value problem (23):

$$\begin{cases} (\mathcal{L} + \lambda \Gamma(\mathbf{x}))F(\mathbf{x}; \lambda) = 0 & \mathbf{x} \in D \\ F(\mathbf{x}; \lambda) = G(\mathbf{x}) & \mathbf{x} \in D^c \end{cases} \quad (\text{A13})$$

830 where we have inserted the additional dependence of F on λ in order to lead directly to the recursive
 831 formulas (20) and (26).

832 *b. Dynkin's formula and finite lag time*

833 We have presented (29) as a mathematically concise approximation to the generator. In practice,
 834 we achieve better numerical stability integrating the generator (A1) to a finite lag time Δt , following
 835 Strahan et al. (2021). The theorem that allows this is called Dynkin's formula (e.g., Oksendal 2003),
 836 which states that for any suitable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and a stopping time θ (not to be confused
 837 with CV coordinates),

$$\mathbb{E}_{\mathbf{x}}[f(\mathbf{X}(\theta))] = f(\mathbf{x}) + \mathbb{E}_{\mathbf{x}} \left[\int_0^\theta \mathcal{L}f(\mathbf{X}(t)) dt \right]. \quad (\text{A14})$$

838 The left-hand side, $\mathbb{E}_{\mathbf{x}}[f(\mathbf{X}(\theta))]$, is known as the *transition operator* $\mathcal{T}^\theta f(\mathbf{x})$, a finite-time version
 839 of the generator. Note that this is a deterministic operator despite θ being a random variable,
 840 because by definition \mathcal{T}^θ only has θ inside of expectations. We can apply Dynkin's formula
 841 to (A13) *before* numerical approximation, setting $\theta = \min(\Delta t, \tau)$. That is, the short trajectory
 842 $\{\mathbf{X}(t) : 0 \leq t \leq \Delta t = 20 \text{ days}\}$ is stopped early if it exits the domain D before Δt . Applying
 843 Dynkin's formula to $F(\mathbf{x}; \lambda)$, we find

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[F(\mathbf{X}(\theta))] &= F(\mathbf{x}) + \mathbb{E}_{\mathbf{x}} \left[\int_0^\theta \mathcal{L}F(\mathbf{X}(t)) dt \right] \\ &= F(\mathbf{x}) - \lambda \mathbb{E}_{\mathbf{x}} \left[\int_0^\theta \Gamma(\mathbf{X}(t)) F(\mathbf{X}(t)) dt \right] \\ \mathcal{T}^\theta F(\mathbf{x}) &= F(\mathbf{x}) - \lambda \mathcal{I}^\theta[\Gamma F](\mathbf{x}) \end{aligned} \quad (\text{A15})$$

844 where \mathcal{I}^θ is shorthand notation for the integral operator on the right. Equation (A15), along with
 845 the boundary conditions $F|_{D^c} = G|_{D^c}$, gives us a linear equation for $F(\mathbf{x})$ that can be solved by
 846 DGA. As outlined in Section 5, we write $F = \hat{F} + f$, where \hat{F} obeys the boundary conditions and

847 f obeys

$$\begin{aligned}
 (\mathcal{T}^\theta - 1)f(\mathbf{x}) + \lambda \mathcal{I}^\theta [\Gamma f](\mathbf{x}) = & \quad (\text{A16}) \\
 - (\mathcal{T}^\theta - 1)\hat{F}(\mathbf{x}) - \lambda \mathcal{I}^\theta [\Gamma \hat{F}](\mathbf{x})
 \end{aligned}$$

848 We then expand $f = \sum_{j=1}^M \xi_j \phi_j(\mathbf{x})$ with basis functions $\{\phi_j\}$ that are zero on D^c , and take μ -
 849 weighted inner products with ϕ_i on both sides to obtain

$$\begin{aligned}
 \sum_{j=1}^M \xi_j \left(\langle \phi_i, (\mathcal{T}^\theta - 1)\phi_j \rangle_\mu + \lambda \langle \phi_i, \mathcal{I}^\theta [\Gamma \phi_j] \rangle_\mu \right) = & \\
 - \langle \phi_i, (\mathcal{T}^\theta - 1)\hat{F} \rangle_\mu - \lambda \langle \phi_i, \mathcal{I}^\theta [\Gamma \hat{F}] \rangle_\mu & \quad (\text{A17})
 \end{aligned}$$

850 Finally, the inner products can be estimated with short trajectories using (30). For two functions
 851 ϕ and ψ , the first left-hand side inner product is approximately

$$\langle \phi, (\mathcal{T}^\theta - 1)\psi \rangle_\mu \approx \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{X}_n) [\psi(\mathbf{X}_n(\theta_n)) - \psi(\mathbf{X}_n)] \quad (\text{A18})$$

852 where θ_n is the sampled first-exit time of the n th trajectory, or Δt if it never exits. The second
 853 left-hand side inner product is approximately

$$\begin{aligned}
 \langle \phi, \mathcal{I}^\theta [\Gamma \psi] \rangle_\mu \approx & \\
 \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{X}_n) \int_0^{\theta_n} \Gamma(\mathbf{X}_n(t)) \psi(\mathbf{X}_n(t)) dt & \quad (\text{A19})
 \end{aligned}$$

854 where the time integral on the right is computed with the trapezoid rule on trajectory, which is
 855 sampled every 0.5 days.

856 Given a fixed Γ and G , and with the inner products in hand, we now have (A17) as a family of
 857 matrix equations with λ a continuous parameter:

$$(P + \lambda Q)\boldsymbol{\xi}(\lambda) = \mathbf{v} + \lambda \mathbf{r}. \quad (\text{A20})$$

858 We can then differentiate in λ and evaluate at $\lambda = 0$ to obtain a ready-to-solve discretization of the
 859 recursion (26):

$$P\xi(0) = \mathbf{v} \quad (\text{A21})$$

$$P\xi'(0) = \mathbf{r} - Q\xi(0) \quad (\text{A22})$$

$$P\xi^{(k)}(0) = -kQ\xi^{(k-1)}(0) \text{ for } k \geq 2 \quad (\text{A23})$$

860 where the k 'th derivative $\xi^{(k)}(0)$ is the coefficient expansion in the basis $\{\phi_j\}$ of the k 'th moment
 861 from (22):

$$\partial_\lambda^k F(\mathbf{x}; 0) = \mathbb{E}_{\mathbf{x}} \left[G(\mathbf{X}(\tau)) \left(\lambda \int_0^\tau \Gamma(\mathbf{X}(s)) ds \right)^k \right] \quad (\text{A24})$$

862 c. Change of measure

863 We now specify how to compute the change of measure from μ (the sampling distribution) to
 864 π (the steady-state distribution), using an adjoint version of the Feynman-Kac formula. Each of
 865 the basis functions ϕ_i has an expectation at time zero with respect to the steady state distribution:
 866 $\mathbb{E}_{\mathbf{X}(0) \sim \pi} [\phi_i(\mathbf{X}(0))] = \int \phi_i(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}$. Evolving the dynamics from 0 to Δt induces another expect-
 867 ation: $\mathbb{E}_{\mathbf{X}(0) \sim \pi} [\phi_i(\mathbf{X}(\Delta t))] = \int \mathcal{T}^{\Delta t} \phi_i(\mathbf{x}) \pi(d\mathbf{x})$. π is the *invariant* distribution, which means that
 868 these two integrals are equal:

$$\int (\mathcal{T}^{\Delta t} - 1) \phi_i(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x} = 0. \quad (\text{A25})$$

869 Furthermore, with a change of measure they can be rewritten with respect to the sampling measure
 870 μ instead of π , so

$$\int (\mathcal{T}^{\Delta t} - 1) \phi_i(\mathbf{x}) \frac{d\pi}{d\mu}(\mathbf{x}) \mu(\mathbf{x}) d\mathbf{x} = 0 \quad (\text{A26})$$

871 The change of measure $\frac{d\pi}{d\mu}(\mathbf{x})$, which we abbreviate $w(\mathbf{x})$, is yet another unknown function which
 872 we expand in the basis as $w(\mathbf{x}) = \sum_j \xi_j \phi_j(\mathbf{x})$. Putting this into the integral and using Monte Carlo,

873 we cast the coefficients ξ_j as the solution to a null eigenvector problem:

$$0 = \int (\mathcal{T}^{\Delta t} - 1)\phi_i(\mathbf{x}) \sum_{j=1}^M c_j(w)\phi_j(\mathbf{x})\mu(d\mathbf{x}) \quad (\text{A27})$$

$$\approx \sum_{j=1}^M c_j(w) \sum_{n=1}^N [\phi_i(\mathbf{X}_n(\Delta t)) - \phi_i(\mathbf{X}_n)]\phi_j(\mathbf{X}_n) \quad (\text{A28})$$

874 This last equation is simply the Fokker-Planck equation, $\mathcal{L}^*\pi = 0$, in weak form and integrated in
 875 time using Dynkin's formula. Note that the matrix elements in (A28) are the transpose of those
 876 in (A18).

877 *d. DGA details*

878 We will provide more details here on our particular construction of basis functions. The partition
 879 $\{S_1, \dots, S_M\}$ to build the basis function library $\phi_j(\mathbf{x}) = \mathbb{1}_{S_j}(\mathbf{x})$, $n = 1, \dots, N$ should be chosen with
 880 a number of considerations in mind. The partition elements should be small enough to accurately
 881 represent the functions they are used to approximate, but large enough to contain sufficient data
 882 to robustly estimate transition probabilities. We form these sets by a hierarchical modification of
 883 K -means clustering on the initial points $\{\mathbf{X}_n\}_{n=1}^N$. K -means is a robust method that can incorporate
 884 new samples by simply identifying the closest centroid, and is commonly used in molecular
 885 dynamics (Pande et al. 2010). However, straightforward application of K -means, as implemented
 886 in the `scikit-learn` software (Pedregosa et al. 2011), can produce a very imbalanced cluster
 887 size distribution, even with empty clusters. This leads to unwanted singularities in the constructed
 888 Markov matrix. To avoid this problem we cluster hierarchically, starting with a coarse clustering
 889 of all points and iteratively refining the larger clusters, at every stage enforcing a minimum cluster
 890 size of five points, until we have the desired number of clusters (M). After clustering on the initial
 891 points $\{\mathbf{X}_n\}$, the other points $\{\mathbf{X}_n(t), 0 < t \leq \Delta t\}$ are placed into clusters using an address tree
 892 produced by the K -means cluster hierarchy. For boundary value problems with a domain D and

893 boundary D^c , we need only cluster points in D , since the basis should be homogeneous. The total
894 number of clusters should scale with data set. In our main results with $N = 5 \times 10^5$, we found
895 $M = 1500$ to be enough basis functions to resolve some of the finer details in the structure of
896 the forecast functions, but not so many as to require an unmanageably deep address tree, which
897 manifests in dramatic slowdown past a certain threshold. At this point, the cluster number is still a
898 manually tuned hyperparameter.

899 Because the committor and lead time obey Dirichlet boundary conditions on $A \cup B$, the basis
900 functions used to construct them should be zero on $A \cup B$, meaning only data points $\mathbf{X}_n \notin A \cup B$
901 should be used to produce the clusters. On the other hand, the steady state distribution has no
902 boundary condition to satisfy, only a global normalization condition. Hence, the basis for the
903 change of measure w must be different from the basis for q^+ and η^+ , with its clusters including all
904 data points in $A \cup B$. Furthermore, the basis must be chosen so that the matrix $\langle (\mathcal{T}^{\Delta t} - 1)\phi_i, \phi_j \rangle$
905 has a nontrivial null space; this is guaranteed by the indicator basis set we use, but can otherwise
906 be guaranteed by including a constant function in the basis.

907 The use of an indicator basis follows the Markov State Modeling literature (Chodera et al. 2006;
908 Pande et al. 2010, e.g.), which has the advantage of simplicity and robustness. In particular, the
909 discretization of $\mathcal{T}^\theta - 1$ is a properly normalized stochastic matrix (with nonnegative entries and
910 rows summing to 1), which guarantees the maximum principle $0 \leq q^+(\mathbf{x}) \leq 1$ and $0 \leq w(\mathbf{x})$ for all
911 data points \mathbf{x} . However, alternative basis sets have been shown to be promising, perhaps with much
912 less data. Thiede et al. (2019) used diffusion maps, while Strahan et al. (2021) used a PCA-like
913 procedure to construct the basis. More generally, there is no requirement to use a linear Galerkin
914 method to solve the Feynman-Kac formulae. More flexible functional forms may have an important
915 role to play as well. In the low-data regime, some preliminary experiments have suggested that
916 Gaussian process regression (GPR) is a useful way to constrain the committor estimate with a prior,

917 following the framework in Bilonis (2016) to solve PDEs with Gaussian processes. As mentioned
918 in the conclusion, there is rapidly growing interest in the use of artificial neural networks to solve
919 PDEs. As with many novel methods, however, DGA is likely to work best on new applications
920 when its simplest form is applied first. This will be our approach in coming experiments on more
921 complex models.

922 **References**

- 923 Berner, J., G. J. Shutts, M. Leutbecher, and T. N. Palmer, 2009: A spectral stochastic kinetic
924 energy backscatter scheme and its impact on flow-dependent predictability in the ecmwf ensemble
925 prediction system. *Journal of the Atmospheric Sciences*, **66 (3)**, 603 – 626, doi:10.1175/
926 2008JAS2677.1, URL <https://journals.ametsoc.org/view/journals/atsc/66/3/2008jas2677.1.xml>.
- 927 Berry, T., J. R. Cressman, Z. Gregurić-Ferenček, and T. Sauer, 2013: Time-scale separation
928 from diffusion-mapped delay coordinates. *SIAM Journal on Applied Dynamical Systems*, **12 (2)**,
929 618–649, doi:10.1137/12088183X, URL <https://doi.org/10.1137/12088183X>, [https://doi.org/](https://doi.org/10.1137/12088183X)
930 [10.1137/12088183X](https://doi.org/10.1137/12088183X).
- 931 Berry, T., D. Giannakis, and J. Harlim, 2015: Nonparametric forecasting of low-dimensional
932 dynamical systems. *Phys. Rev. E*, **91**, 032 915, doi:10.1103/PhysRevE.91.032915.
- 933 Bilonis, I., 2016: Probabilistic solvers for partial differential equations. *arXiv: Probability*.
- 934 Binzel, R. P., 2000: The torino impact hazard scale. *Planetary and Space Science*,
935 **48 (4)**, 297–303, doi:[https://doi.org/10.1016/S0032-0633\(00\)00006-4](https://doi.org/10.1016/S0032-0633(00)00006-4), URL [https://www.](https://www.sciencedirect.com/science/article/pii/S0032063300000064)
936 [sciencedirect.com/science/article/pii/S0032063300000064](https://www.sciencedirect.com/science/article/pii/S0032063300000064).
- 937 Birner, T., and P. D. Williams, 2008: Sudden stratospheric warmings as noise-induced transitions.
938 *Journal of the Atmospheric Sciences*, **65 (10)**, 3337–3343, doi:10.1175/2008JAS2770.1.

- 939 Bolton, T., and L. Zanna, 2019: Applications of deep learning to ocean data in-
940 ference and subgrid parameterization. *Journal of Advances in Modeling Earth Sys-*
941 *tems*, **11** (1), 376–399, doi:<https://doi.org/10.1029/2018MS001472>, URL [https://agupubs.](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001472)
942 [onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001472](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001472), [https://agupubs.onlinelibrary.wiley.](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018MS001472)
943 [com/doi/pdf/10.1029/2018MS001472](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018MS001472).
- 944 Bouchet, F., J. Laurie, and O. Zaboronski, 2011: Control and instanton trajectories for ran-
945 dom transitions in turbulent flows. *Journal of Physics: Conference Series*, **318** (2), 022041,
946 doi:10.1088/1742-6596/318/2/022041, URL <https://doi.org/10.1088/1742-6596/318/2/022041>,
947 <https://doi.org/10.1088/1742-6596/318/2/022041>.
- 948 Bouchet, F., J. Laurie, and O. Zaboronski, 2014: Langevin dynamics, large deviations and instan-
949 tons for the quasi-geostrophic model and two-dimensional euler equations. *Journal of Statis-*
950 *tical Physics*, **156**, 1066–1092, doi:10.1007/s10955-014-1052-5, URL [https://doi.org/10.1007/](https://doi.org/10.1007/s10955-014-1052-5)
951 [s10955-014-1052-5](https://doi.org/10.1007/s10955-014-1052-5).
- 952 Bouchet, F., J. Rolland, and E. Simonnet, 2019a: Rare event algorithm links transitions in
953 turbulent flows with activated nucleations. *Physical Review Letters*, **122** (7), 074502, doi:
954 [10.1103/PhysRevLett.122.074502](https://doi.org/10.1103/PhysRevLett.122.074502).
- 955 Bouchet, F., J. Rolland, and J. Wouters, 2019b: Rare event sampling methods. *Chaos: An Inter-*
956 *disciplinary Journal of Nonlinear Science*, **29** (8), 080402, doi:10.1063/1.5120509.
- 957 Bowman, G. R., V. S. Pande, and F. Noé, 2013: *An introduction to Markov state models and*
958 *their application to long timescale molecular simulation*, Vol. 797. Springer Science & Business
959 Media.

960 Carleo, G., and M. Troyer, 2017: Solving the quantum many-body problem with arti-
961 cial neural networks. *Science*, **355** (6325), 602–606, doi:10.1126/science.aag2302, URL
962 <https://science.sciencemag.org/content/355/6325/602>, [https://science.sciencemag.org/content/](https://science.sciencemag.org/content/355/6325/602.full.pdf)
963 [355/6325/602.full.pdf](https://science.sciencemag.org/content/355/6325/602.full.pdf).

964 Charlton, A. J., and L. M. Polvani, 2007: A new look at stratospheric sudden warmings. part
965 i: Climatology and modeling benchmarks. *Journal of Climate*, **20** (3), 449–469, doi:10.1175/
966 JCLI3996.1.

967 Charney, J. G., and J. G. DeVore, 1979: Multiple Flow Equilibria in the At-
968 mosphere and Blocking. *Journal of the Atmospheric Sciences*, **36** (7), 1205–1216,
969 doi:10.1175/1520-0469(1979)036<1205:MFEITA>2.0.CO;2, URL [https://doi.org/10.1175/](https://doi.org/10.1175/1520-0469(1979)036<1205:MFEITA>2.0.CO;2)
970 [1520-0469\(1979\)036<1205:MFEITA>2.0.CO;2](https://doi.org/10.1175/1520-0469(1979)036<1205:MFEITA>2.0.CO;2), [https://journals.ametsoc.org/jas/article-pdf/](https://journals.ametsoc.org/jas/article-pdf/36/7/1205/3420739/1520-0469(1979)036%5C_1205%5C_mfeita%5C_2%5C_0%5C_co%5C_2.pdf)
971 [36/7/1205/3420739/1520-0469\(1979\)036%5C_1205%5C_mfeita%5C_2%5C_0%5C_co%5C_2.pdf](https://journals.ametsoc.org/jas/article-pdf/36/7/1205/3420739/1520-0469(1979)036%5C_1205%5C_mfeita%5C_2%5C_0%5C_co%5C_2.pdf).

972 Chattopadhyay, A., M. Mustafa, P. Hassanzadeh, E. Bach, and K. Kashinath, 2021: Towards physi-
973 cally consistent data-driven weather forecasting: Integrating data assimilation with equivariance-
974 preserving spatial transformers in a case study with era5. *Geoscientific Model Development Dis-*
975 *cussions*, **2021**, 1–23, doi:10.5194/gmd-2021-71, URL [https://gmd.copernicus.org/preprints/](https://gmd.copernicus.org/preprints/gmd-2021-71/)
976 [gmd-2021-71/](https://gmd.copernicus.org/preprints/gmd-2021-71/).

977 Chattopadhyay, A., E. Nabizadeh, and P. Hassanzadeh, 2020: Analog forecast-
978 ing of extreme-causing weather patterns using deep learning. *Journal of Ad-*
979 *vances in Modeling Earth Systems*, **12** (2), e2019MS001958, doi:[https://doi.](https://doi.org/10.1029/2019MS001958)
980 [org/10.1029/2019MS001958](https://doi.org/10.1029/2019MS001958), URL [https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001958)
981 [2019MS001958](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001958), e2019MS001958 10.1029/2019MS001958, [https://agupubs.onlinelibrary.](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2019MS001958)
982 [wiley.com/doi/pdf/10.1029/2019MS001958](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2019MS001958).

983 Chen, N., D. Giannakis, R. Herbei, and A. J. Majda, 2014: An mcmc algorithm for parameter esti-
984 mation in signals with hidden intermittent instability. *SIAM/ASA Journal on Uncertainty Quan-*
985 *tification*, **2 (1)**, 647–669, doi:10.1137/130944977, URL <https://doi.org/10.1137/130944977>,
986 <https://doi.org/10.1137/130944977>.

987 Chen, N., and A. J. Majda, 2020: Predicting observed and hidden extreme events in complex
988 nonlinear dynamical systems with partial observations and short training time series. *Chaos: An*
989 *Interdisciplinary Journal of Nonlinear Science*, **30 (3)**, 033–101, doi:10.1063/1.5122199, URL
990 <https://doi.org/10.1063/1.5122199>, <https://doi.org/10.1063/1.5122199>.

991 Chodera, J. D., and F. Noé, 2014: Markov state models of biomolecular conformational dynamics.
992 *Current Opinion in Structural Biology*, **25**, 135 – 144, doi:[https://doi.org/10.1016/j.sbi.2014.04.](https://doi.org/10.1016/j.sbi.2014.04.002)
993 [002](http://www.sciencedirect.com/science/article/pii/S0959440X14000426), URL <http://www.sciencedirect.com/science/article/pii/S0959440X14000426>, theory and
994 simulation / Macromolecular machines.

995 Chodera, J. D., W. C. Swope, J. W. Pitera, and K. A. Dill, 2006: Long-time protein folding
996 dynamics from short-time molecular dynamics simulations. *Multiscale Modeling & Simulation*,
997 **5 (4)**, 1214–1226, doi:10.1137/06065146X.

998 Christiansen, B., 2000: Chaos, quasiperiodicity, and interannual variability: Studies of a strato-
999 spheric vacillation model. *Journal of the Atmospheric Sciences*, **57 (18)**, 3161–3173, doi:
1000 [10.1175/1520-0469\(2000\)057<3161:CQAIVS>2.0.CO;2](https://doi.org/10.1175/1520-0469(2000)057<3161:CQAIVS>2.0.CO;2).

1001 Crommelin, D. T., 2003: Regime transitions and heteroclinic connections in a barotropic
1002 atmosphere. *Journal of the Atmospheric Sciences*, **60 (2)**, 229 – 246, doi:10.
1003 [1175/1520-0469\(2003\)060<0229:RTAHCI>2.0.CO;2](https://doi.org/10.1175/1520-0469(2003)060<0229:RTAHCI>2.0.CO;2), URL [https://journals.ametsoc.org/view/](https://journals.ametsoc.org/view/journals/atsc/60/2/1520-0469_2003_060_0229_rtahci_2.0.co_2.xml)
1004 [journals/atsc/60/2/1520-0469_2003_060_0229_rtahci_2.0.co_2.xml](https://journals.ametsoc.org/view/journals/atsc/60/2/1520-0469_2003_060_0229_rtahci_2.0.co_2.xml).

1005 DelSole, T., and B. F. Farrell, 1995: A stochastically excited linear system as a model for quasi-
1006 geostrophic turbulence: Analytic results for one- and two-layer fluids. *Journal of the Atmospheric*
1007 *Sciences*, **52 (14)**, 2531–2547, doi:10.1175/1520-0469(1995)052<2531:ASELSA>2.0.CO;2.

1008 Dematteis, G., T. Grafke, M. Onorato, and E. Vanden-Eijnden, 2019: Experimental evidence
1009 of hydrodynamic instantons: The universal route to rogue waves. *Phys. Rev. X*, **9**, 041 057,
1010 doi:10.1103/PhysRevX.9.041057, URL <https://link.aps.org/doi/10.1103/PhysRevX.9.041057>.

1011 Dematteis, G., T. Grafke, and E. Vanden-Eijnden, 2018: Rogue waves and large deviations in deep
1012 sea. *Proceedings of the National Academy of Sciences*, **115 (5)**, 855–860, doi:10.1073/pnas.
1013 1710670115.

1014 Durrett, R., 2013: *Probability: Theory and Examples*. Cambridge University Press.

1015 E, W., T. Li, and E. Vanden-Eijnden, 2019: *Applied stochastic analysis*, Vol. 199. American
1016 Mathematical Soc.

1017 E., W., and E. Vanden-Eijnden, 2006: Towards a Theory of Transition Paths. *Journal of Sta-*
1018 *tistical Physics*, **123 (3)**, 503, doi:10.1007/s10955-005-9003-9, URL [https://doi.org/10.1007/](https://doi.org/10.1007/s10955-005-9003-9)
1019 [s10955-005-9003-9](https://doi.org/10.1007/s10955-005-9003-9).

1020 Esler, J. G., and M. Mester, 2019: Noise-induced vortex-splitting stratospheric sudden warmings.
1021 *Quarterly Journal of the Royal Meteorological Society*, **145 (719)**, 476–494, doi:[https://doi.org/](https://doi.org/10.1002/qj.3443)
1022 [10.1002/qj.3443](https://doi.org/10.1002/qj.3443), URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3443>, [https://](https://rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.3443)
1023 rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.3443.

1024 Farazmand, M., and T. P. Sapsis, 2017: A variational approach to probing extreme events in turbu-
1025 lent dynamical systems. *Science Advances*, **3 (9)**, doi:10.1126/sciadv.1701533, URL <https://>

1026 advances.sciencemag.org/content/3/9/e1701533, [https://advances.sciencemag.org/content/3/9/](https://advances.sciencemag.org/content/3/9/e1701533.full.pdf)
1027 e1701533.full.pdf.

1028 Finkel, J., D. S. Abbot, and J. Weare, 2020: Path Properties of Atmospheric Transitions: Illustration
1029 with a Low-Order Sudden Stratospheric Warming Model. *Journal of the Atmospheric*
1030 *Sciences*, **77** (7), 2327–2347, doi:10.1175/JAS-D-19-0278.1, URL [https://doi.org/10.1175/](https://doi.org/10.1175/JAS-D-19-0278.1)
1031 [JAS-D-19-0278.1](https://doi.org/10.1175/JAS-D-19-0278.1), [https://journals.ametsoc.org/jas/article-pdf/77/7/2327/4958190/jasd190278.](https://journals.ametsoc.org/jas/article-pdf/77/7/2327/4958190/jasd190278.pdf)
1032 pdf.

1033 Fitzsimmons, P., and J. Pitman, 1999: Kac’s moment formula and the feynman–kac for-
1034 mula for additive functionals of a markov process. *Stochastic Processes and their Appli-*
1035 *cations*, **79** (1), 117–134, doi:[https://doi.org/10.1016/S0304-4149\(98\)00081-7](https://doi.org/10.1016/S0304-4149(98)00081-7), URL <https://www.sciencedirect.com/science/article/pii/S0304414998000817>.
1036

1037 Frank, N., and S. Fischer, 2008: Transition networks for modeling the kinetics of conformational
1038 change in macromolecules. *Current Opinion in Structural Biology*, **18**, 154–163, doi:10.1016/
1039 j.sbi.2008.01.008.

1040 Franzke, C., and A. J. Majda, 2006: Low-order stochastic mode reduction for a prototype atmo-
1041 spheric gcm. *Journal of the Atmospheric Sciences*, **63** (2), 457–479, doi:10.1175/JAS3633.1.

1042 Giannakis, D., A. Kolchinskaya, D. Krasnov, and J. Schumacher, 2018: Koopman analysis of the
1043 long-term evolution in a turbulent convection cell. *Journal of Fluid Mechanics*, **847**, 735–767,
1044 doi:10.1017/jfm.2018.297.

1045 Giannakis, D., and A. J. Majda, 2012: Nonlinear laplacian spectral analysis for time series with
1046 intermittency and low-frequency variability. *Proceedings of the National Academy of Sciences*,

1047 **109 (7)**, 2222–2227, doi:10.1073/pnas.1118984109, <https://www.pnas.org/content/109/7/2222>.
1048 full.pdf.

1049 Gottwald, G. A., D. T. Crommelin, and C. L. E. Franzke, 2016: Stochastic climate theory.
1050 1612.07474.

1051 Han, J., A. Jentzen, and W. E., 2018: Solving high-dimensional partial differential equations
1052 using deep learning. *Proceedings of the National Academy of Sciences*, **115 (34)**, 8505–8510,
1053 doi:10.1073/pnas.1718942115, URL <https://www.pnas.org/content/115/34/8505>, <https://www.pnas.org/content/115/34/8505.full.pdf>.

1055 Hasselmann, K., 1976: Stochastic climate models part i. theory. *Tellus*, **28 (6)**, 473–485, doi:
1056 10.3402/tellusa.v28i6.11316.

1057 Helfmann, L., J. Heitzig, P. Koltai, J. Kurths, and C. Schütte, 2021: Statistical analysis of tipping
1058 pathways in agent-based models. 2103.02883.

1059 Hoffman, R. N., J. M. Henderson, S. M. Leidner, C. Grassotti, and T. Nehr Korn, 2006: The response
1060 of damaging winds of a simulated tropical cyclone to finite-amplitude perturbations of different
1061 variables. *Journal of the Atmospheric Sciences*, **63 (7)**, 1924 – 1937, doi:10.1175/JAS3720.1,
1062 URL <https://journals.ametsoc.org/view/journals/atsc/63/7/jas3720.1.xml>.

1063 Holton, J. R., and C. Mass, 1976: Stratospheric vacillation cycles. *Journal of the Atmospheric*
1064 *Sciences*, **33 (11)**, 2218–2225, doi:10.1175/1520-0469(1976)033<2218:SVC>2.0.CO;2.

1065 Karatzas, I., and S. E. Shreve, 1998: *Brownian Motion and Stochastic Calculus*. Springer.

1066 Kashinath, K., and Coauthors, 2021: Physics-informed machine learning: case studies for weather
1067 and climate modelling. *Philosophical Transactions of the Royal Society A: Mathematical,*

1068 *Physical and Engineering Sciences*, **379 (2194)**, 20200 093, doi:10.1098/rsta.2020.0093, URL
1069 <https://royalsocietypublishing.org/doi/abs/10.1098/rsta.2020.0093>.

1070 Khoo, Y., J. Lu, and L. Ying, 2018: Solving for high-dimensional committor functions
1071 using artificial neural networks. *Research in the Mathematical Sciences*, **6**, doi:10.1007/
1072 s40687-018-0160-2, URL <https://doi.org/10.1007/s40687-018-0160-2>.

1073 Legras, B., and M. Ghil, 1985: Persistent anomalies, blocking and variations in at-
1074 mospheric predictability. *Journal of Atmospheric Sciences*, **42 (5)**, 433 – 471, doi:10.
1075 1175/1520-0469(1985)042<0433:PABAVI>2.0.CO;2, URL [https://journals.ametsoc.org/view/
1076 journals/atsc/42/5/1520-0469_1985_042_0433_pabavi_2_0_co_2.xml](https://journals.ametsoc.org/view/journals/atsc/42/5/1520-0469_1985_042_0433_pabavi_2_0_co_2.xml).

1077 Li, H., Y. Khoo, Y. Ren, and L. Ying, 2020: Solving for high dimensional committor functions
1078 using neural network with online approximation to derivatives. 2012.06727.

1079 Li, Q., B. Lin, and W. Ren, 2019: Computing committor functions for the study of rare events using
1080 deep learning. *The Journal of Chemical Physics*, **151 (5)**, 054 112, doi:10.1063/1.5110439, URL
1081 <https://doi.org/10.1063/1.5110439>.

1082 Lin, K. K., and F. Lu, 2021: Data-driven model reduction, wiener projections, and the
1083 koopman-mori-zwanzig formalism. *Journal of Computational Physics*, **424**, 109 864, doi:
1084 <https://doi.org/10.1016/j.jcp.2020.109864>, URL [https://www.sciencedirect.com/science/article/
1085 pii/S0021999120306380](https://www.sciencedirect.com/science/article/pii/S0021999120306380).

1086 Lorenz, E. N., 1963: Deterministic nonperiodic flow. *Journal of Atmospheric Sciences*, **20 (2)**, 130
1087 – 141, doi:10.1175/1520-0469(1963)020<0130:DNF>2.0.CO;2, URL [https://journals.ametsoc.
1088 org/view/journals/atsc/20/2/1520-0469_1963](https://journals.ametsoc.org/view/journals/atsc/20/2/1520-0469_1963).

- 1089 Lorpaiboon, C., E. H. Thiede, R. J. Webber, J. Weare, and A. R. Dinner, 2020: Integrated
1090 variational approach to conformational dynamics: A robust strategy for identifying eigenfunc-
1091 tions of dynamical operators. *The Journal of Physical Chemistry B*, **124** (42), 9354–9364,
1092 doi:10.1021/acs.jpcc.0c06477, URL <https://doi.org/10.1021/acs.jpcc.0c06477>.
- 1093 Lucarini, V., and A. Gritsun, 2020: A new mathematical framework for atmospheric blocking
1094 events. *Climate Dynamics*, **54** (1), 575–598.
- 1095 Lucente, D., S. Duffner, C. Herbert, J. Rolland, and F. Bouchet, 2019: Machine learning of
1096 committor functions for predicting high impact climate events. *Climate Informatics*, Paris,
1097 France, URL <https://hal.archives-ouvertes.fr/hal-02322370>.
- 1098 Maiocchi, C. C., V. Lucarini, A. Gritsun, and G. Pavliotis, 2020: Unstable Periodic Orbits Sam-
1099 pling in Climate Models. *EGU General Assembly Conference Abstracts*, 18823, EGU General
1100 Assembly Conference Abstracts.
- 1101 Majda, A. J., and D. Qi, 2018: Strategies for reduced-order models for predicting the statistical
1102 responses and uncertainty quantification in complex turbulent dynamical systems. *SIAM Review*,
1103 **60** (3), 491–549, doi:10.1137/16M1104664, URL <https://doi.org/10.1137/16M1104664>, <https://doi.org/10.1137/16M1104664>.
- 1105 Majda, A. J., I. Timofeyev, and E. Vanden Eijnden, 2001: A mathematical framework for
1106 stochastic climate models. *Communications on Pure and Applied Mathematics*, **54** (8),
1107 891–974, doi:<https://doi.org/10.1002/cpa.1014>, URL [https://onlinelibrary.wiley.com/doi/abs/](https://onlinelibrary.wiley.com/doi/abs/10.1002/cpa.1014)
1108 [10.1002/cpa.1014](https://onlinelibrary.wiley.com/doi/pdf/10.1002/cpa.1014), <https://onlinelibrary.wiley.com/doi/pdf/10.1002/cpa.1014>.
- 1109 Mardt, A., L. Pasuali, H. Wu, and F. Noé, 2018: Vampnets for deep learning of molecular kinetics.
1110 *Nature Communications*, **9**, doi:10.1038/s41467-017-02388-1, URL <https://doi.org/10.1038/>

1111 s41467-017-02388-1.

1112 Matsuno, T., 1971: A Dynamical Model of the Stratospheric Sudden Warming. *Journal*
1113 *of the Atmospheric Sciences*, **28 (8)**, 1479–1494, doi:10.1175/1520-0469(1971)028<1479:
1114 ADMOTS>2.0.CO;2, URL [https://doi.org/10.1175/1520-0469\(1971\)028<1479:ADMOTS>2.](https://doi.org/10.1175/1520-0469(1971)028<1479:ADMOTS>2.0.CO;2)
1115 [0.CO;2](https://doi.org/10.1175/1520-0469(1971)028<1479:ADMOTS>2.0.CO;2), [https://journals.ametsoc.org/jas/article-pdf/28/8/1479/3417422/1520-0469\(1971\)028\](https://journals.ametsoc.org/jas/article-pdf/28/8/1479/3417422/1520-0469(1971)028%5B1479%5D_admots%5B2%5D_co%5B2%5D.pdf)
1116 [_1479\admots\2\0_co\2.pdf](https://journals.ametsoc.org/jas/article-pdf/28/8/1479/3417422/1520-0469(1971)028%5B1479%5D_admots%5B2%5D_co%5B2%5D.pdf).

1117 Metzner, P., C. Schutte, and E. Vanden-Eijnden, 2006: Illustration of transition path theory
1118 on a collection of simple examples. *The Journal of Chemical Physics*, **125 (8)**, 1–2, doi:
1119 10.1063/1.2335447.

1120 Metzner, P., C. Schutte, and E. Vanden-Eijnden, 2009: Transition path theory for markov jump
1121 processes. *Multiscale Modeling and Simulation*, **7 (3)**, 1192–1219, doi:10.1137/070699500.

1122 Miron, P., F. Beron-Vera, L. Helfmann, and P. Koltai, 2021: Transition paths of marine debris and
1123 the stability of the garbage patches. *Chaos: An Interdisciplinary Journal of Nonlinear Science*,
1124 accepted for publication.

1125 Mohamad, M. A., and T. P. Sapsis, 2018: Sequential sampling strategy for extreme event
1126 statistics in nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*,
1127 **115 (44)**, 11 138–11 143, doi:10.1073/pnas.1813263115, URL [https://www.pnas.org/content/](https://www.pnas.org/content/115/44/11138)
1128 [115/44/11138](https://www.pnas.org/content/115/44/11138), <https://www.pnas.org/content/115/44/11138.full.pdf>.

1129 Ngwira, C. M., and Coauthors, 2013: Simulation of the 23 july 2012 extreme
1130 space weather event: What if this extremely rare cme was earth directed? *Space*
1131 *Weather*, **11 (12)**, 671–679, doi:<https://doi.org/10.1002/2013SW000990>, URL <https://agupubs>.

1132 onlinelibrary.wiley.com/doi/abs/10.1002/2013SW000990, <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1002/2013SW000990>.

1134 Noé, F., C. Schütte, E. Vanden-Eijnden, L. Reich, and T. R. Weigl, 2009: Constructing the equilib-
1135 rium ensemble of folding pathways from short off-equilibrium simulations. *Proceedings of the*
1136 *National Academy of Sciences*, **106 (45)**, 19 011–19 016, doi:10.1073/pnas.0905466106, URL
1137 <https://www.pnas.org/content/106/45/19011>, <https://www.pnas.org/content/106/45/19011.full.pdf>.

1139 Noé, F., and C. Clementi, 2017: Collective variables for the study of long-time kinetics from molec-
1140 ular trajectories: theory and methods. *Current Opinion in Structural Biology*, **43**, 141–147, doi:
1141 <https://doi.org/10.1016/j.sbi.2017.02.006>, URL <https://www.sciencedirect.com/science/article/pii/S0959440X17300301>, theory and simulation • Macromolecular assemblies.

1143 Oksendal, B., 2003: *Stochastic Differential Equations: An Introduction with Applications*.
1144 Springer.

1145 Pande, V. S., K. Beauchamp, and G. R. Bowman, 2010: Everything you wanted to know about
1146 markov state models but were afraid to ask. *Methods*, **52 (1)**, 99–105, URL <https://doi.org/10.1016/j.ymeth.2010.06.002>.

1148 Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine learning in Python. *Journal of Machine*
1149 *Learning Research*, **12**, 2825–2830.

1150 Plotkin, D. A., R. J. Webber, M. E. O’Neill, J. Weare, and D. S. Abbot, 2019: Maximizing
1151 simulated tropical cyclone intensity with action minimization. *Journal of Advances in Modeling*
1152 *Earth Systems*, **11 (4)**, 863–891, doi:10.1029/2018MS001419.

1153 Porta Mana, P., and L. Zanna, 2014: Toward a stochastic parameterization of ocean mesoscale
1154 eddies. *Ocean Modelling*, **79**, 1–20, doi:<https://doi.org/10.1016/j.ocemod.2014.04.002>, URL
1155 <https://www.sciencedirect.com/science/article/pii/S1463500314000420>.

1156 Ragone, F., and F. Bouchet, 2020: Computation of extreme values of time averaged observables
1157 in climate models with large deviation techniques. *Journal of Statistical Physics*, **179** (5), 1637–
1158 1665, doi:10.1007/s10955-019-02429-7, URL <https://doi.org/10.1007/s10955-019-02429-7>.

1159 Ragone, F., J. Wouters, and F. Bouchet, 2018: Computation of extreme heat waves in climate
1160 models using a large deviation algorithm. *Proceedings of the National Academy of Sciences*,
1161 **115** (1), 24–29, doi:10.1073/pnas.1712645115, <https://www.pnas.org/content/115/1/24.full.pdf>.

1162 Raissi, M., P. Perdikaris, and G. Karniadakis, 2019: Physics-informed neural networks: A deep
1163 learning framework for solving forward and inverse problems involving nonlinear partial differ-
1164 ential equations. *Journal of Computational Physics*, **378**, 686–707, doi:<https://doi.org/10.1016/j.jcp.2018.10.045>, URL <https://www.sciencedirect.com/science/article/pii/S0021999118307125>.

1166 Ruzmaikin, A., J. Lawrence, and C. Cadavid, 2003: A simple model of stratospheric dynamics
1167 including solar variability. *Journal of Climate*, **16**, 1593–1600, doi:10.1175/2007JCLI2119.1.

1168 Sabeerali, C. T., R. S. Ajayamohan, D. Giannakis, and A. J. Majda, 2017: Extraction and prediction
1169 of indices for monsoon intraseasonal oscillations: an approach based on nonlinear laplacian
1170 spectral analysis. *Climate Dynamics*, **49** (9), 3031–3050, doi:10.1007/s00382-016-3491-y.

1171 Sapsis, T. P., 2021: Statistics of extreme events in fluid flows and waves. *Annual Review of Fluid*
1172 *Mechanics*, **53** (1), 85–111, doi:10.1146/annurev-fluid-030420-032810, URL [https://doi.org/10.](https://doi.org/10.1146/annurev-fluid-030420-032810)
1173 [1146/annurev-fluid-030420-032810](https://doi.org/10.1146/annurev-fluid-030420-032810), <https://doi.org/10.1146/annurev-fluid-030420-032810>.

1174 Schaller, N., J. Sillmann, J. Anstey, E. M. Fischer, C. M. Grams, and S. Russo, 2018: Influence of
1175 blocking on northern european and western russian heatwaves in large climate model ensembles.
1176 *Environmental Research Letters*, **13** (5), 054 015, doi:10.1088/1748-9326/aaba55, URL https:
1177 //doi.org/10.1088%2F1748-9326%2Faaba55.

1178 Simonnet, E., J. Rolland, and F. Bouchet, 2020: Multistability and rare spontaneous transitions be-
1179 tween climate and jet configurations in a barotropic model of the jovian mid-latitude troposphere.
1180 2009.09913.

1181 Sjoberg, J. P., and T. Birner, 2014: Stratospheric wave–mean flow feedbacks and sudden
1182 stratospheric warmings in a simple model forced by upward wave activity flux. *Journal*
1183 *of the Atmospheric Sciences*, **71** (11), 4055 – 4071, doi:10.1175/JAS-D-14-0113.1, URL
1184 https://journals.ametsoc.org/view/journals/atsc/71/11/jas-d-14-0113.1.xml.

1185 Strahan, J., A. Antoszewski, C. Lorpaiboon, B. P. Vani, J. Weare, and A. R. Dinner, 2021:
1186 Long-time-scale predictions from short-trajectory data: A benchmark analysis of the trp-cage
1187 miniprotein. *Journal of Chemical Theory and Computation*, **17** (5), 2948–2963, doi:10.1021/
1188 acs.jctc.0c00933, URL https://doi.org/10.1021/acs.jctc.0c00933, pMID: 33908762, https://doi.
1189 org/10.1021/acs.jctc.0c00933.

1190 Tantet, A., F. R. van der Burgt, and H. A. Dijkstra, 2015: An early warning indicator for atmo-
1191 spheric blocking events using transfer operators. *Chaos: An Interdisciplinary Journal of Nonlin-*
1192 *ear Science*, **25** (3), 036 406, doi:10.1063/1.4908174, URL https://doi.org/10.1063/1.4908174,
1193 https://doi.org/10.1063/1.4908174.

1194 Thiede, E., D. Giannakis, A. R. Dinner, and J. Weare, 2019: Approximation of dynamical quantities
1195 using trajectory data. *arXiv:1810.01841 [physics.data-an]*, 1–24, doi:1810.01841.

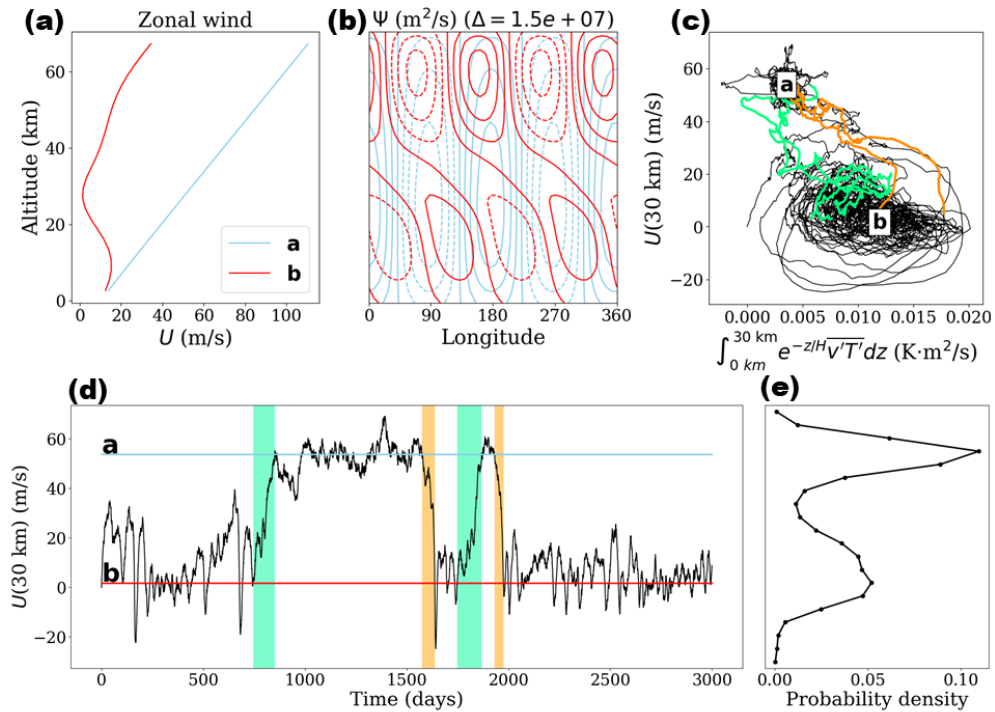
- 1196 Tibshirani, R., 1996: Regression shrinkage and selection via the lasso. *Journal of the Royal*
1197 *Statistical Society: Series B (Methodological)*, **58** (1), 267–288, doi:[https://doi.org/10.](https://doi.org/10.1111/j.2517-6161.1996.tb02080.x)
1198 [10.1111/j.2517-6161.1996.tb02080.x](https://doi.org/10.1111/j.2517-6161.1996.tb02080.x), URL [https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/](https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/j.2517-6161.1996.tb02080.x)
1199 [j.2517-6161.1996.tb02080.x](https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1996.tb02080.x), [https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.](https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1996.tb02080.x)
1200 [1996.tb02080.x](https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1996.tb02080.x).
- 1201 Timmermann, A., F.-F. Jin, and J. Abshagen, 2003: A nonlinear theory for el niño bursting. *Jour-*
1202 *nal of the Atmospheric Sciences*, **60** (1), 152 – 165, doi:[10.1175/1520-0469\(2003\)060<0152:](https://doi.org/10.1175/1520-0469(2003)060<0152:ANTFEN>2.0.CO;2)
1203 [ANTFEN>2.0.CO;2](https://doi.org/10.1175/1520-0469(2003)060<0152:ANTFEN>2.0.CO;2), URL [https://journals.ametsoc.org/view/journals/atsc/60/1/1520-0469_](https://journals.ametsoc.org/view/journals/atsc/60/1/1520-0469_2003)
1204 [2003](https://journals.ametsoc.org/view/journals/atsc/60/1/1520-0469_2003).
- 1205 Vanden-Eijnden, E., and W. E, 2010: Transition-path theory and path-finding algorithms for the
1206 study of rare events. *Annual Review of Physical Chemistry*, **61** (1), 391–420, doi:[10.1146/](https://doi.org/10.1146/annurev.physchem.040808.090412)
1207 [annurev.physchem.040808.090412](https://doi.org/10.1146/annurev.physchem.040808.090412).
- 1208 Vanden-Eijnden, E., and J. Weare, 2013: Data assimilation in the low noise regime with application
1209 to the kuroshio. *Monthly Weather Review*, **141** (6), 1822–1841, doi:[10.1175/MWR-D-12-00060.](https://doi.org/10.1175/MWR-D-12-00060.1)
1210 [1](https://doi.org/10.1175/MWR-D-12-00060.1).
- 1211 Vitart, F., and A. W. Robertson, 2018: The sub-seasonal to seasonal prediction project (s2s)
1212 and the prediction of extreme events. *npj Climate and Atmospheric Science*, **1**, URL [https:](https://doi.org/10.1038/s41612-018-0013-0)
1213 [//doi.org/10.1038/s41612-018-0013-0](https://doi.org/10.1038/s41612-018-0013-0).
- 1214 Wan, Z. Y., P. Vlachas, P. Koumoutsakos, and T. Sapsis, 2018: Data-assisted reduced-order
1215 modeling of extreme events in complex dynamical systems. *PLOS ONE*, **13** (5), 1–22, doi:
1216 [10.1371/journal.pone.0197704](https://doi.org/10.1371/journal.pone.0197704), URL <https://doi.org/10.1371/journal.pone.0197704>.

- 1217 Weare, J., 2009: Particle filtering with path sampling and an application to a bimodal ocean current
1218 model. *Journal of Computational Physics*, **228** (12), 4312 – 4331, doi:[https://doi.org/10.1016/j.](https://doi.org/10.1016/j.jcp.2009.02.033)
1219 [jcp.2009.02.033](https://doi.org/10.1016/j.jcp.2009.02.033).
- 1220 Webber, R. J., D. A. Plotkin, M. E. O’Neill, D. S. Abbot, and J. Weare, 2019: Practical rare event
1221 sampling for extreme mesoscale weather. *Chaos*, **29** (5), 053 109, doi:10.1063/1.5081461.
- 1222 Yasuda, Y., F. Bouchet, and A. Venaille, 2017: A new interpretation of vortex-split sudden
1223 stratospheric warmings in terms of equilibrium statistical mechanics. *Journal of the Atmospheric*
1224 *Sciences*, **74** (12), 3915–3936, doi:10.1175/JAS-D-17-0045.1.
- 1225 Yoden, S., 1987a: Bifurcation properties of a stratospheric vacillation model. *Journal of the At-*
1226 *mospheric Sciences*, **44** (13), 1723–1733, doi:10.1175/1520-0469(1987)044<1723:BPOASV>
1227 2.0.CO;2.
- 1228 Yoden, S., 1987b: Dynamical Aspects of Stratospheric Vacillations in a Highly
1229 Truncated Model. *Journal of the Atmospheric Sciences*, **44** (24), 3683–3695,
1230 doi:10.1175/1520-0469(1987)044<3683:DAOSVI>2.0.CO;2, URL [https://doi.org/10.1175/](https://doi.org/10.1175/1520-0469(1987)044<3683:DAOSVI>2.0.CO;2)
1231 [1520-0469\(1987\)044<3683:DAOSVI>2.0.CO;2](https://doi.org/10.1175/1520-0469(1987)044<3683:DAOSVI>2.0.CO;2).
- 1232 Zhang, F., and J. A. Sippel, 2009: Effects of moist convection on hurricane predictability. *Journal*
1233 *of the Atmospheric Sciences*, **66** (7), 1944 – 1961, doi:10.1175/2009JAS2824.1, URL [https:](https://journals.ametsoc.org/view/journals/atsc/66/7/2009jas2824.1.xml)
1234 [//journals.ametsoc.org/view/journals/atsc/66/7/2009jas2824.1.xml](https://journals.ametsoc.org/view/journals/atsc/66/7/2009jas2824.1.xml).
- 1235 Zwanzig, R., 2001: *Nonequilibrium statistical mechanics*. Oxford University Press.

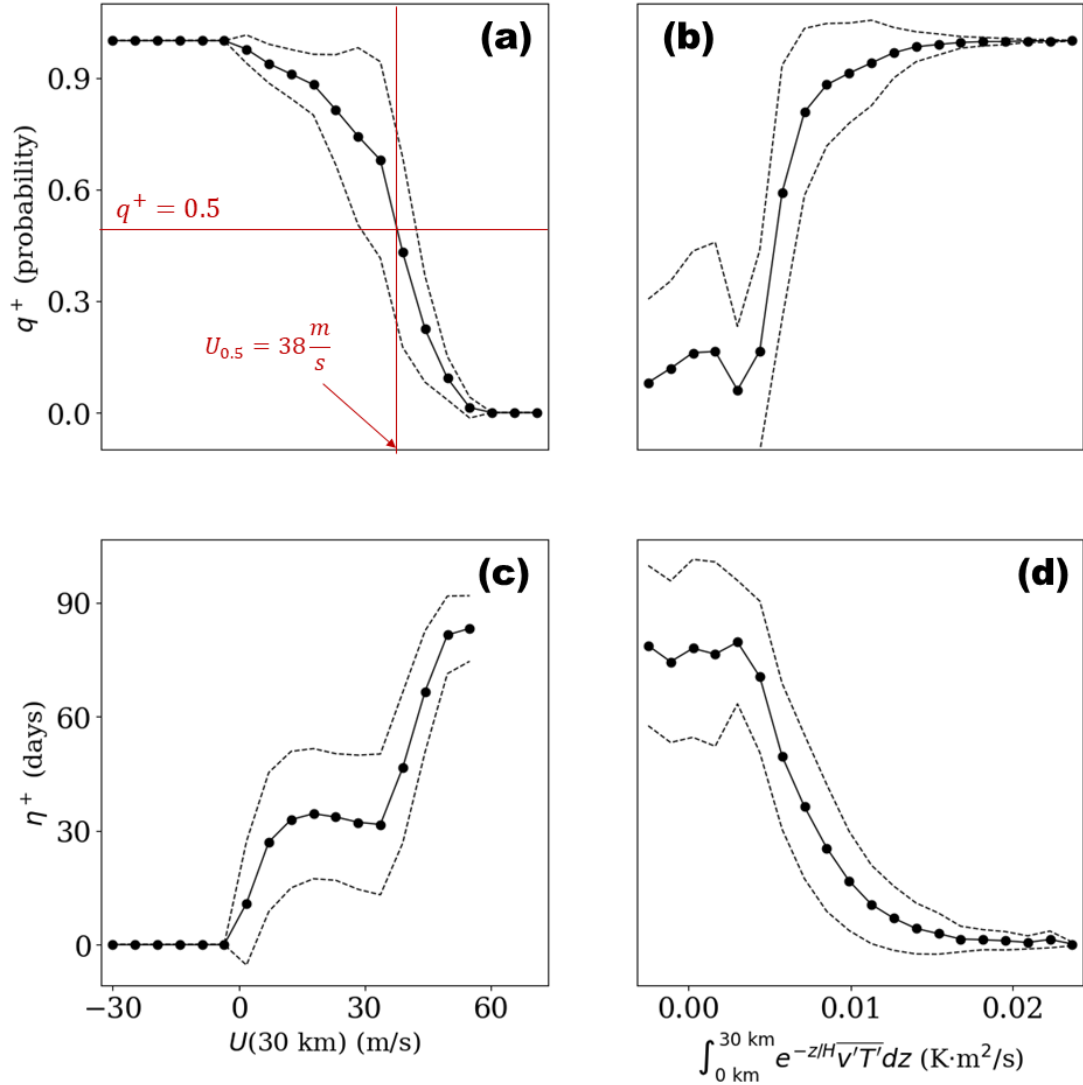
LIST OF FIGURES

1236		
1237	Fig. 1.	Illustration of the two stable states of the Holton-Mass model and transitions between them. (a) Zonal wind profiles of the radiatively maintained strong vortex (the fixed point a , blue) which increases linearly with altitude, and the weak vortex (the fixed point b , red) which dips close to zero in the mid-stratosphere. (b) Streamfunction contours are overlaid for the two equilibria a and b . (c) Parametric plot of a control simulation in a 2-dimensional state space projection, including two transitions from <i>A</i> to <i>B</i> (orange) and <i>B</i> to <i>A</i> (green). (d) Time series of $U(30\text{ km})$ from the same simulation. (e) The steady state density projected onto $U(30\text{ km})$ 63
1238		
1239		
1240		
1241		
1242		
1243		
1244		
1245	Fig. 2.	One-dimensional projections of the forward committor (first row) and lead time to <i>B</i> (second row). These functions depend on all $d = 75$ degrees of freedom in the model, but we have averaged across $d - 1 = 74$ dimensions to visualize them as rough functions of two single degrees of freedom: $U(30\text{ km})$ (first column) and integrated heat flux up to 30 km, IHF (second column). Panel (a) additionally marks the $q^+ = \frac{1}{2}$ threshold and the corresponding value of zonal wind. 64
1246		
1247		
1248		
1249		
1250		
1251	Fig. 3.	The density, committor, and lead time as functions of zonal wind and integrated heat flux. Panel (a) projects the steady state distribution $\pi(\mathbf{x})$ onto the two-dimensional subspace (U, IHF) at 30 km. The white regions surrounding the gray are unphysical states with negligible probability. Panels (b) and (c) display the committor and lead time in the same space. A horizontal transect marks the level $U(30\text{ km}) = 38.5\text{ m/s}$, where q^+ according to U only is 0.5. Panels (d) and (e) show ensembles initialized from two points θ_0 and θ_1 along the transect, verifying that their committor and lead time values differ from their values according to U , in a way that is predictable due to considering IHF in addition to U 65
1252		
1253		
1254		
1255		
1256		
1257		
1258		
1259	Fig. 4.	Committor and lead time as independent coordinates. This figure inverts the functions in Figure 3, considering the zonal wind and integrated heat flux as functions of committor and lead time. The two-dimensional space they span is the essential goal of forecasting. Panel (a) shows the steady state distribution on this subspace, which is peaked near a and b (darker shading), weaker in the "bridge" region between them, and completely negligible the white regions unexplored by data. Panels (b) and (c) display zonal wind and heat flux in color as functions of the committor and lead time. 66
1260		
1261		
1262		
1263		
1264		
1265		
1266	Fig. 5.	Projection of the forward committor onto a large collection of altitude-dependent physical variables. The top left panel shows heatmaps of q^+ as a function of U and z ; white regions denote where $U(z)$ is negligibly observed. The top middle panel shows the standard deviation in q^+ as a function of U and z ; this uncertainty stems from the remaining 74 model dimensions. The right-hand panel displays the total mean-squared error due to the projection for each altitude, i.e., $\sqrt{S[f; \theta]}$ from Equation (14). A low value indicates that this level is ideal for prediction. The following rows show the same quantities for other physical variables: streamfunction magnitude, eddy enstrophy, background PV gradient, eddy PV flux, and LASSO. 67
1267		
1268		
1269		
1270		
1271		
1272		
1273		
1274		
1275	Fig. 6.	Results of LASSO regression of the forward committor with linear and nonlinear input features. Panel (a) shows the coefficients when q^+ is regressed as a function of only the variables at a given altitude, and panel (b) shows the corresponding correlation score. 21.5 km seems the most predictive (where $z \equiv 0$ at the tropopause, not the surface). Panel (c) shows the coefficient structure when all altitudes are considered simultaneously. Most of the nonzero coefficients appear between 15-22 km, distinguishing that range as highly relevant for prediction. 68
1276		
1277		
1278		
1279		
1280		
1281		

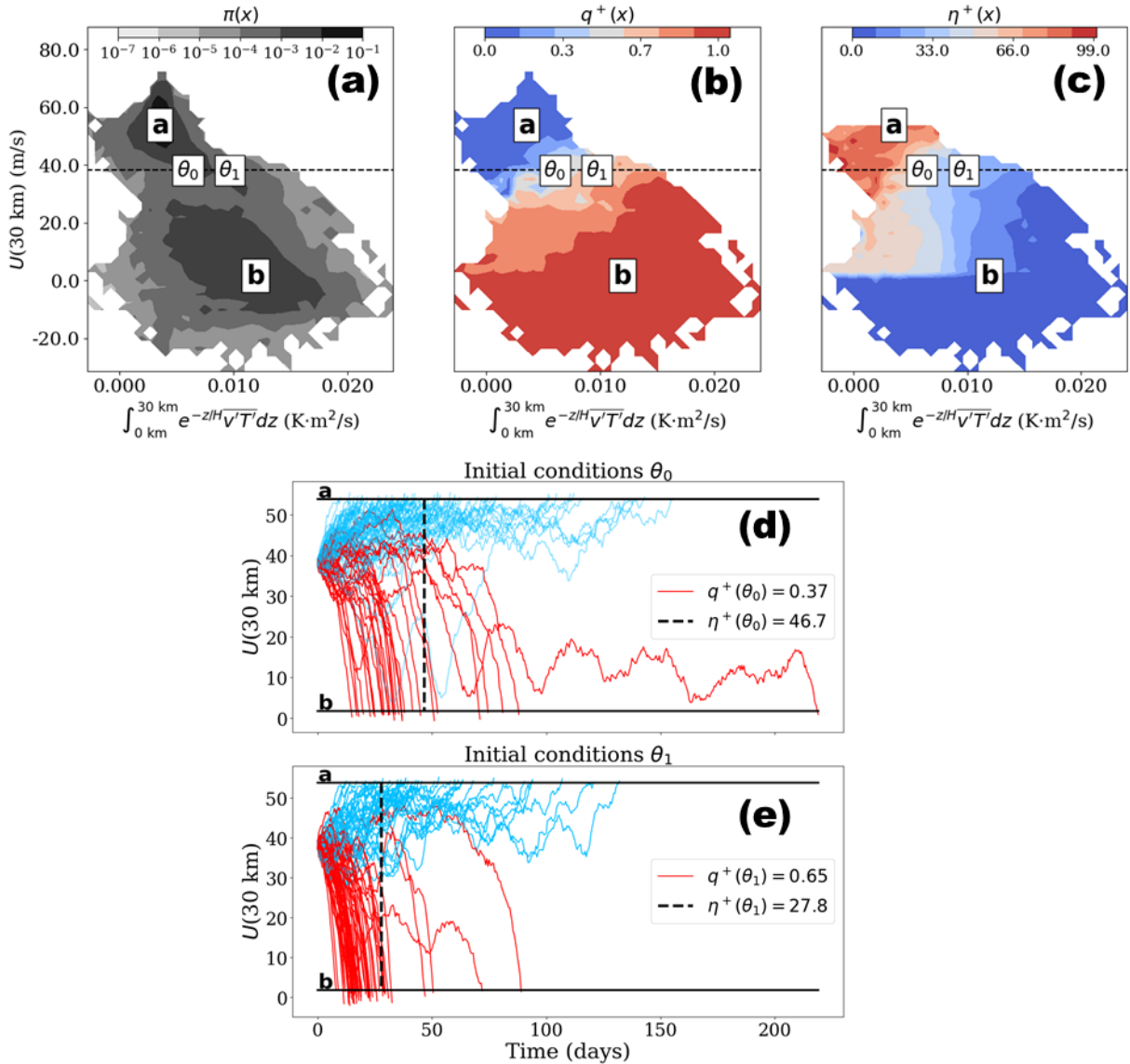
1282 **Fig. 7. Fidelity of DGA.** For several DGA parameter values of N (the number of data points), M
1283 (the number of basis functions) and lag time, we plot the committor calculated from DGA
1284 and from the long control simulation, both as a function of $U(30 \text{ km})$. The mean-square
1285 difference ε in the legend is used as a global error estimate for DGA. 69



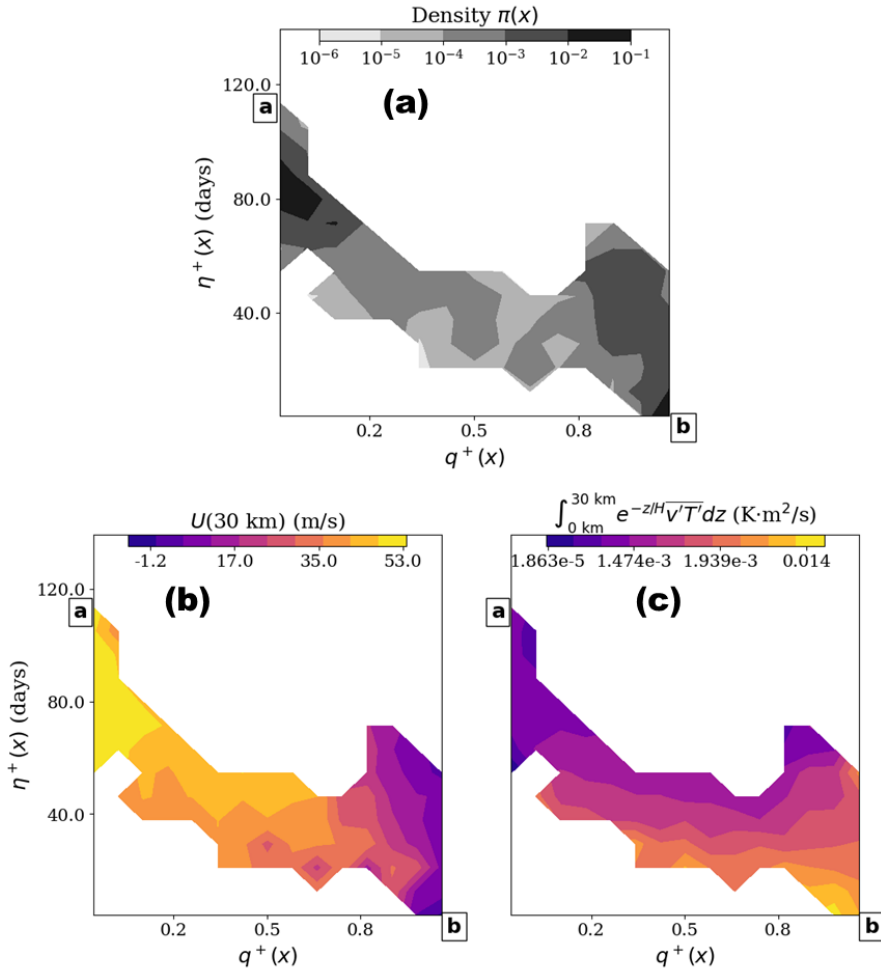
1286 FIG. 1. **Illustration of the two stable states of the Holton-Mass model and transitions between them.** (a)
 1287 Zonal wind profiles of the radiatively maintained strong vortex (the fixed point **a**, blue) which increases linearly
 1288 with altitude, and the weak vortex (the fixed point **b**, red) which dips close to zero in the mid-stratosphere. (b)
 1289 Streamfunction contours are overlaid for the two equilibria **a** and **b**. (c) Parametric plot of a control simulation
 1290 in a 2-dimensional state space projection, including two transitions from *A* to *B* (orange) and *B* to *A* (green). (d)
 1291 Time series of $U(30 \text{ km})$ from the same simulation. (e) The steady state density projected onto $U(30 \text{ km})$.



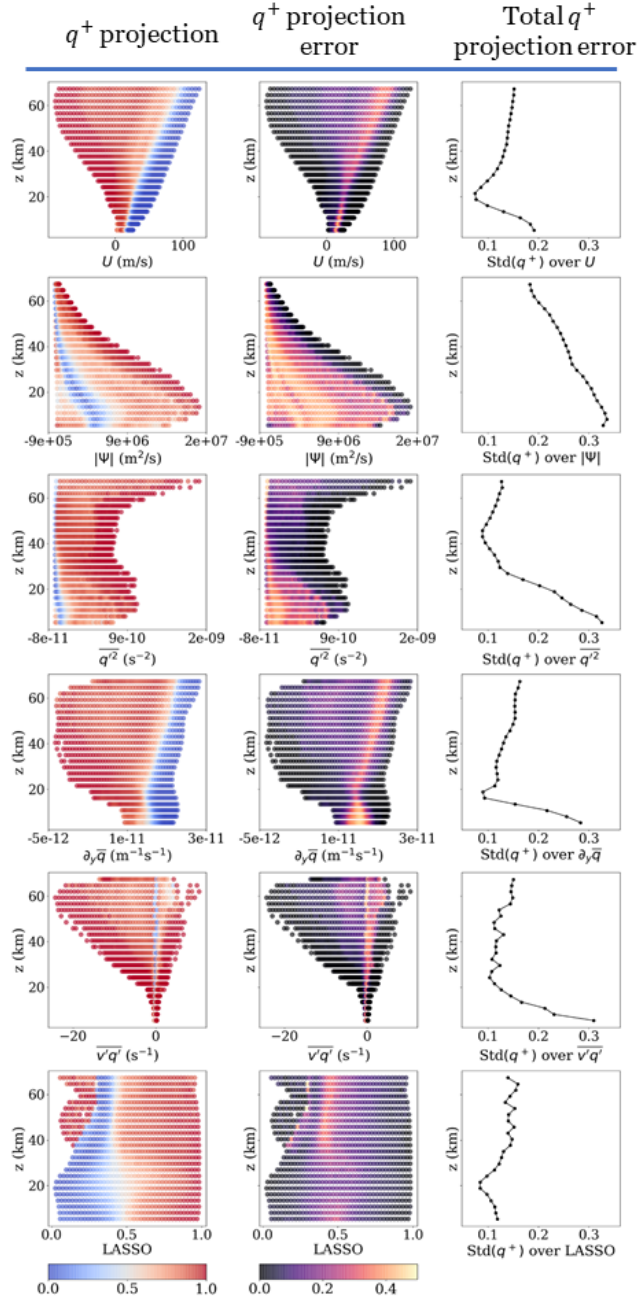
1292 FIG. 2. **One-dimensional projections of the forward committor (first row) and lead time to B (second**
 1293 **row).** These functions depend on all $d = 75$ degrees of freedom in the model, but we have averaged across
 1294 $d - 1 = 74$ dimensions to visualize them as rough functions of two single degrees of freedom: $U(30 \text{ km})$ (first
 1295 column) and integrated heat flux up to 30 km, IHF (second column). Panel (a) additionally marks the $q^+ = \frac{1}{2}$
 1296 threshold and the corresponding value of zonal wind.



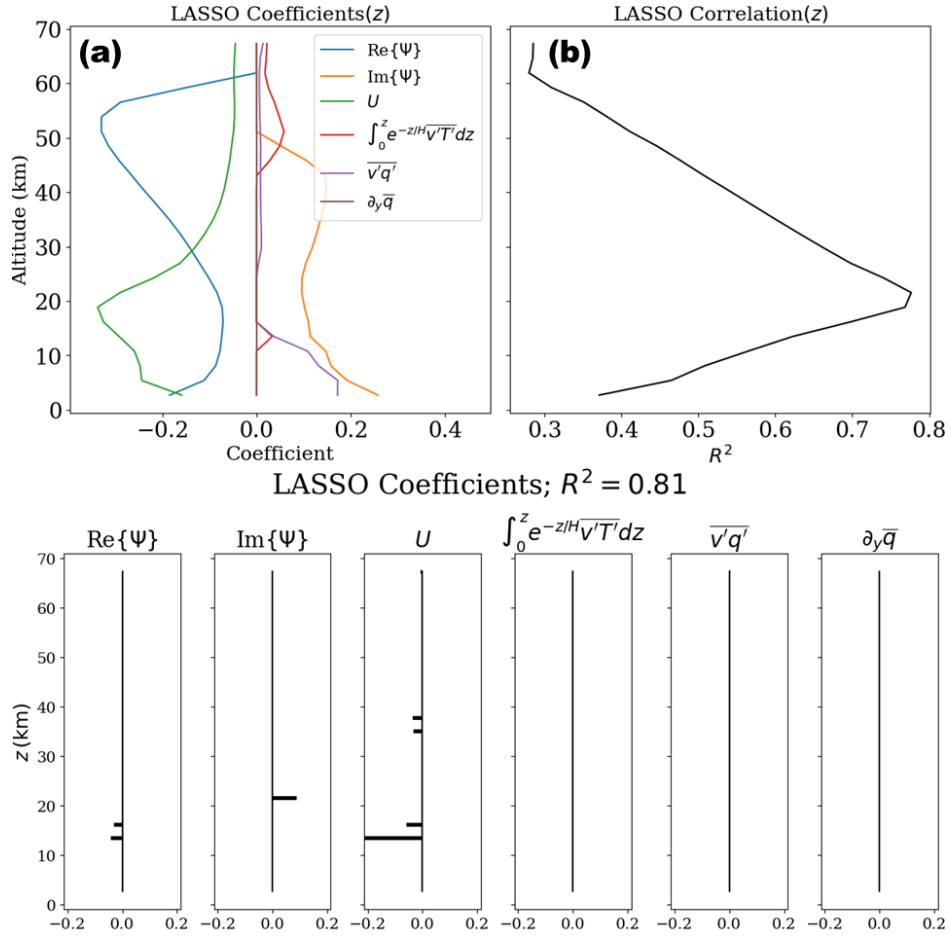
1297 **FIG. 3. The density, committor, and lead time as functions of zonal wind and integrated heat flux.** Panel
 1298 (a) projects the steady state distribution $\pi(\mathbf{x})$ onto the two-dimensional subspace (U, IHF) at 30 km. The white
 1299 regions surrounding the gray are unphysical states with negligible probability. Panels (b) and (c) display the
 1300 committor and lead time in the same space. A horizontal transect marks the level $U(30 \text{ km}) = 38.5 \text{ m/s}$, where
 1301 q^+ according to U only is 0.5. Panels (d) and (e) show ensembles initialized from two points θ_0 and θ_1 along
 1302 the transect, verifying that their committor and lead time values differ from their values according to U , in a way
 1303 that is predictable due to considering IHF in addition to U .



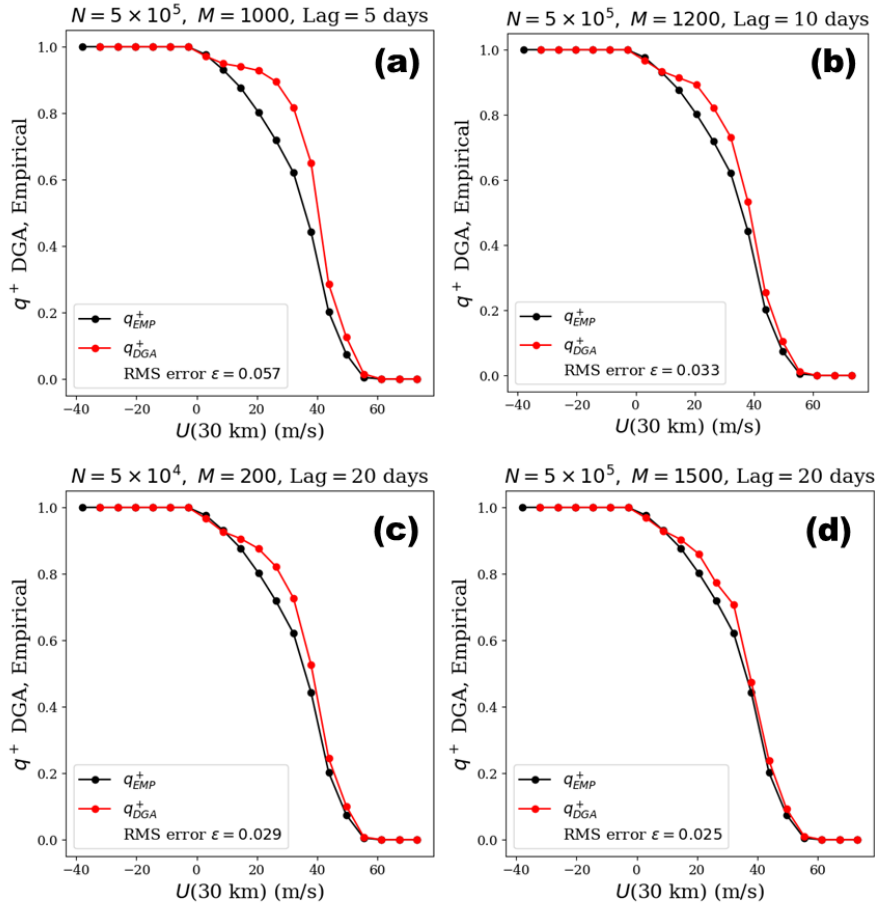
1304 **FIG. 4. Committor and lead time as independent coordinates.** This figure inverts the functions in Figure 3,
 1305 considering the zonal wind and integrated heat flux as functions of committor and lead time. The two-dimensional
 1306 space they span is the essential goal of forecasting. Panel (a) shows the steady state distribution on this subspace,
 1307 which is peaked near **a** and **b** (darker shading), weaker in the "bridge" region between them, and completely
 1308 negligible the white regions unexplored by data. Panels (b) and (c) display zonal wind and heat flux in color as
 1309 functions of the committor and lead time.



1310 FIG. 5. **Projection of the forward committor onto a large collection of altitude-dependent physical**
 1311 **variables.** The top left panel shows heatmaps of q^+ as a function of U and z ; white regions denote where
 1312 $U(z)$ is negligibly observed. The top middle panel shows the standard deviation in q^+ as a function of U and
 1313 z ; this uncertainty stems from the remaining 74 model dimensions. The right-hand panel displays the total
 1314 mean-squared error due to the projection for each altitude, i.e., $\sqrt{S[f; \theta]}$ from Equation (14). A low value
 1315 indicates that this level is ideal for prediction. The following rows show the same quantities for other physical
 1316 variables: streamfunction magnitude, eddy enstrophy, background PV gradient, eddy PV flux, and LASSO.



1317 **FIG. 6. Results of LASSO regression of the forward committor with linear and nonlinear input features.**
 1318 Panel (a) shows the coefficients when q^+ is regressed as a function of only the variables at a given altitude,
 1319 and panel (b) shows the corresponding correlation score. 21.5 km seems the most predictive (where $z \equiv 0$ at
 1320 the tropopause, not the surface). Panel (c) shows the coefficient structure when all altitudes are considered
 1321 simultaneously. Most of the nonzero coefficients appear between 15-22 km, distinguishing that range as highly
 1322 relevant for prediction.



1323 FIG. 7. **Fidelity of DGA.** For several DGA parameter values of N (the number of data points), M (the number of
1324 basis functions) and lag time, we plot the committor calculated from DGA and from the long control simulation,
1325 both as a function of $U(30 \text{ km})$. The mean-square difference ϵ in the legend is used as a global error estimate for
1326 DGA.