

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](mailto:jfinkel@uchicago.edu)

## ABSTRACT

13 Nonlinear atmospheric dynamics produce rare events that are hard to predict and attribute due to  
14 many interacting degrees of freedom. Sudden stratospheric warming event is a model example.  
15 Approximately once every other year, the winter polar vortex in the boreal stratosphere rapidly  
16 breaks down, inducing a shift in midlatitude surface weather patterns persisting for up to 2-  
17 3 months. In principle, lengthy numerical simulations can be used to predict and understand  
18 these rare transitions. For complex models, however, the cost of the direct numerical simulation  
19 approach is often prohibitive. We describe an alternative approach which only requires relatively  
20 short-duration computer simulations of the system. The methodology is illustrated by applying  
21 it to a prototype model of an SSW event developed by Holton and Mass (1976) and driven with  
22 stochastic forcing. While highly idealized, the model captures the essential nonlinear dynamics of  
23 SSWs and exhibits the key forecasting challenge: the dramatic separation in timescales between  
24 the dynamics of a single event and the return time between successive events. We compute optimal  
25 forecasts of sudden warming events and quantify the limits of predictability. Statistical analysis  
26 relates these optimal forecasts to a small number of interpretable physical variables. Remarkably,  
27 we are able to estimate these quantities using a data set of simulations much shorter than the  
28 timescale of the warming event. This methodology is designed to take full advantage of the high-  
29 dimensional data from models and observations, and can be employed to find detailed predictors  
30 of many complex rare events arising in climate dynamics.

## 31 **1. Introduction**

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

43 Large ensemble simulations are the most detailed source of data to assess the frequency, intensity,  
44 and correlates of extreme weather events (e.g., Schaller et al. 2018). A single simulation must  
45 span decades to incorporate the possible impacts of climate change and decadal-scale variability  
46 and the full state space of a climate model may be billions of dimensions large, depending on grid  
47 resolution. Unfortunately, the data-richness of a long simulation comes at the cost of sample size:  
48 even the largest ensembles are limited to tens or hundreds of members as a matter of computational  
49 necessity. Under stationary background parameters and some ergodicity properties, a single  
50 simulation will eventually sample state space thoroughly and provide all relevant statistics. In  
51 practice, however, simulations are often not run long enough to reach steady state, and furthermore  
52 one may wish to change parameters over time. A much larger number of independent ensemble  
53 members would then be needed to quantify the effects of initial conditions, changing climatology,

54 feedbacks, and unresolved high-frequency variability with statistical confidence (Sillmann et al.  
55 2017; Webber et al. 2019).

56 While the last decade has seen exciting progress in the development of targeted rare event  
57 simulation in geophysical contexts (Hoffman et al. 2006; Weare 2009; Bouchet et al. 2011, 2014;  
58 Vanden-Eijnden and Weare 2013; Chen et al. 2014; Yasuda et al. 2017; Farazmand and Sapsis  
59 2017; Dematteis et al. 2018; Mohamad and Sapsis 2018; Dematteis et al. 2019; Webber et al.  
60 2019; Bouchet et al. 2019a,b; Plotkin et al. 2019; Simonnet et al. 2020; Ragone and Bouchet  
61 2020; Sapsis 2021), predicting long time-scale behavior of complex dynamical systems remains  
62 a difficult task. A traditional approach to addressing this issue is through dimensional reduction  
63 techniques which seek to replace an expensive, high-fidelity model with a lower-dimensional and  
64 less costly model. Physics-based reduced-order models have a long and very successful history in  
65 atmospheric science, especially as prototypes of chaos and multistability (Lorenz 1963; Charney  
66 and DeVore 1979; Legras and Ghil 1985; Crommelin 2003; Timmermann et al. 2003; Ruzmaikin  
67 et al. 2003). Observationally, regime behavior has been diagnosed by projecting the empirical  
68 steady-state distributions onto leading EOFs (e.g., Crommelin 2003). More recently, significant  
69 attention has been paid to data-based dimensional reduction techniques that use data generated  
70 by the high-fidelity model to specify a more quantitatively accurate reduced-order model (e.g.,  
71 Giannakis et al. 2018; Berry et al. 2015; Sabeerali et al. 2017; Majda and Qi 2018; Wan et al.  
72 2018; Bolton and Zanna 2019; Chattopadhyay et al. 2020; Chen and Majda 2020). However the  
73 reduced-order model is derived, it can subsequently be thoroughly interrogated by direct computer  
74 simulation.

75 We advance an alternative computational approach to predicting and understanding rare events  
76 without sacrificing model fidelity. Like data-informed reduced order modeling, our method relies  
77 on data generated by a high-fidelity model. However, unlike dimensional reduction techniques, our

78 approach focuses on computing specific quantities of interest rather than on capturing all aspects of  
79 a very complicated, high dimensional dynamical system. In particular we will compute estimators  
80 of statistically optimal forecasts using a data set of many short forward simulations. To accomplish  
81 this we represent these forecasts as solutions to Feynman-Kac equations. In the continuous time  
82 limit, these become partial differential equations (PDE) with a number of independent variables  
83 equal to the dimension of the model state space. It is therefore hopeless to solve the equations  
84 using any standard spatial discretization. As we demonstrate nonetheless, the equations can be  
85 solved with remarkable accuracy via an expansion in a basis of functions informed by the data set.  
86 Importantly, our approach to solving these equations is independent of the model used to generate  
87 the data, avoiding unrealistic simplifications or structural assumptions.

88 As typical examples of the forecasts computable within our framework, we focus on the  
89 probability that a warming event occurs before a return to a “typical” state, as well as the expected  
90 time that it takes for that event to occur. Both quantities depend on the initial condition, and are  
91 therefore functions over all of state space. We will follow the convention in computational statistical  
92 mechanics and refer to these as the committor function and mean first passage time (MFPT)  
93 respectively. The committor has been computed previously for low dimensional atmospheric  
94 models in Tantet et al. (2015); Lucente et al. (2019); Finkel et al. (2020). Forecasts like these  
95 quantify the risk associated with an event given the current state of the system. They also encode  
96 important information regarding the rare event itself.

97 Even putting aside the difficulty of computing the committor and MFPT, they still must be  
98 ‘decoded’; knowledge of these functions does not automatically reveal insights into the fundamental  
99 causes or precursors of a rare event. Nor are they easily applied to observations of a limited  
100 collection of variables. They are, after all, complicated functions of a high dimensional model  
101 state space. In Section 5 we will demonstrate a detailed statistical analysis of our computed

102 committor function aimed at identifying a relatively small subset of the original variables capable  
103 of describing the committor (in the sense defined below).

104 We illustrate our approach on the highly simplified Holton-Mass model (Holton and Mass 1976;  
105 Christiansen 2000) with stochastic velocity perturbations in the spirit of Birner and Williams  
106 (2008). The Holton-Mass model is well-understood dynamically in light of decades of analysis  
107 and experiments, yet complex enough to present the essential computational difficulties of proba-  
108 bilistic forecasting and test our methods for addressing them. Despite the challenges posed by its  
109 75-dimensional state space, our computational framework can indeed accurately characterize of  
110 extreme events with unprecedented detail using only a data set of short model simulations. In the  
111 future, the same methodology could be applied to query the properties of more complex models  
112 where less theoretical understanding is available.

113 Section 2 reviews the dynamical model, and section 3 describes a general class of methods which  
114 we then apply to our problem specifically in section 4. We present the results in section 5, including  
115 a discussion of optimal forecasting and physical insights gleaned from our approach. We then lay  
116 out future prospects and conclude in section 6.

## 117 **2. Holton-Mass model**

118 Holton and Mass (1976) devised a simple model of the stratosphere aimed at reproducing  
119 observed intra-seasonal oscillations of the polar vortex, which they termed “stratospheric vacil-  
120 lation cycles.” Earlier SSW models, originating with that of Matsuno (1971), proposed upward-  
121 propagating planetary waves as the major source of disturbance to the vortex. This was a significant  
122 step outside the bounds of the nonacceleration theorem of Charney and Drazin (1961), which stated  
123 that the vortex would be robust to disturbances under a variety of conditions. While Matsuno (1971)  
124 used impulsive forcing from the troposphere as the source of planetary waves, Holton and Mass

125 (1976) suggested that stationary tropospheric forcing, if large enough, could lead to an oscillatory  
126 response, purely through dynamics internal to the stratosphere.

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

134 Yoden (1987a) found that for a certain range of parameter settings, these two effects balance each  
135 other to create two distinct stable regimes: a strong vortex with zonal wind close to the radiative  
136 equilibrium profile, and a weak vortex with a possibly oscillatory wind profile. We focus our study  
137 on this bistable setting as a prototypical model of atmospheric regime behavior. The transition  
138 from strong to weak vortex state captures the essential dynamics of an SSW. The methodology  
139 presented here, using only observed short trajectories, can be applied equally to any of these models  
140 as well as observational data, which the reader should keep in mind as we present the specifics of  
141 the present application.

142 The Holton-Mass model takes the linearized quasigeostrophic potential vorticity (QGPV) equa-  
143 tion for a perturbation streamfunction  $\psi'(x, y, z, t)$  on top of a zonal mean flow  $\bar{u}(y, z, t)$ , and  
144 projects these two fields onto a single zonal wavenumber  $k = 2/(a \cos 60^\circ)$  and a single meridional  
145 wavenumber  $\ell = 3/a$ , where  $a$  is the Earth's radius. The resulting ansatz is

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

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

146 which is fully determined by the reduced state space  $U(z, t)$ , and  $\Psi(z, t)$ , the latter being complex.

147 Inserting this into the linearized QGPV equations yields the coupled PDE system

$$\left[ -\left( \mathcal{G}^2(k^2 + \ell^2) + \frac{1}{4} \right) + \frac{\partial^2}{\partial z^2} \right] \frac{\partial \Psi}{\partial t} \quad (2)$$

$$= \left[ \left( \frac{\alpha}{4} - \frac{\alpha_z}{2} - i\mathcal{G}^2 k\beta \right) - \alpha_z \frac{\partial}{\partial z} - \alpha \frac{\partial^2}{\partial z^2} \right] \Psi$$

$$+ \left\{ ik\varepsilon \left[ \left( k^2 \mathcal{G}^2 + \frac{1}{4} \right) - \frac{\partial}{\partial z} + \frac{\partial^2}{\partial z^2} \right] U \right\} \Psi - ik\varepsilon \frac{\partial^2 \Psi}{\partial z^2} U$$

$$\left( -\mathcal{G}^2 \ell^2 - \frac{\partial}{\partial z} + \frac{\partial^2}{\partial z^2} \right) \frac{\partial U}{\partial t} = [(\alpha_z - \alpha)U_z^R - \alpha U_{zz}^R] \quad (3)$$

$$- \left[ (\alpha_z - \alpha) \frac{\partial}{\partial z} + \alpha \frac{\partial^2}{\partial z^2} \right] U + \frac{\varepsilon k \ell^2}{2} e^z \text{Im} \left\{ \Psi \frac{\partial^2 \Psi^*}{\partial z^2} \right\}$$

148 where we have nondimensionalized the equations with the parameter  $\mathcal{G}^2 = H^2 N^2 / (f_0^2 L^2)$  in order

149 to create a homogeneously shaped dataset more suited to our analysis. Boundary conditions are

150 prescribed at the bottom of the stratosphere, which in this model corresponds to  $z = 0$ , and the top

151 of the stratosphere  $z_{top} = 70 \text{ km}$ .

$$\Psi(0, t) = \frac{gh}{f_0} \quad \Psi(z_{top}, t) = 0 \quad (4)$$

$$U(0, t) = U^R(0) \quad \partial_z U(z_{top}, t) = \partial_z U^R(z_{top})$$

152 The vortex-stabilizing influence is represented by  $\alpha(z)$ , the altitude-dependent cooling coefficient,

153 and the linear relaxation profile  $U^R(z) = U^R(0) + \frac{\gamma}{1000}z$ , which forces the vortex toward radiative

154 equilibrium. Here  $\gamma = O(1)$  is the vertical wind shear in m/s/km. The competing force of wave

155 perturbation is encoded through the lower boundary condition  $\Psi(0, t) = gh/f_0$ .

156 Detailed bifurcation analysis of the model by both Yoden (1987a) and Christiansen (2000) in

157  $(\gamma, h)$  space revealed the bifurcations that lead to bistability, vacillations, and ultimately quasiperi-

158 odicity and chaos. Here we will focus on an intermediate parameter setting of  $\gamma = 1.5 \text{ m/s/km}$

159 and  $h = 38.5 \text{ m}$ , where two stable states coexist: a strong vortex with  $U$  closely following  $U^R$

160 and an almost barotropic streamfunction, as well as a weak vortex with  $U$  dipping close to zero

161 at an intermediate altitude and a disturbed streamfunction with strong westward phase tilt. The  
162 two stable equilibria, which we call **a** and **b**, are represented in the first row of Figure 1 by their  
163  $z$ -dependent zonal wind and streamfunction profiles.

164 The two equilibria can be interpreted as two different winter climatologies, one with a strong  
165 vortex and one with a weak vortex susceptible to vacillation cycles. To explore transitions between  
166 these two states, we follow Birner and Williams (2008) and modify the Holton-Mass equations  
167 with small additive noise in the  $U$  variable to mimic momentum perturbations by smaller scale  
168 Rossby waves, gravity waves, and other unresolved sources. While the details of the additive noise  
169 are ad hoc, this approach can be more rigorously justified through the Mori-Zwanzig formalism  
170 (Zwanzig 2001). Because many hidden degrees of freedom are being projected onto the low-  
171 dimensional space of the Holton-Mass model, the dynamics on small observable subspaces can be  
172 considered stochastic. This is the perspective taken in stochastic parameterization of turbulence  
173 and other high-dimensional chaotic systems (Hasselmann 1976; DelSole and Farrell 1995; Franzke  
174 and Majda 2006; Majda et al. 2001; Gottwald et al. 2016). More sophisticated parameterizations  
175 would surely influence the results (Hu et al. 2019), but would not present a fundamental problem  
176 for our purely data-driven numerical method.

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

181 vector

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

182 Boundary terms are not included because they are constrained by the boundary conditions. The  
 183 deterministic system can be written  $d\mathbf{X}(t)/dt = m(\mathbf{X}(t))$  for a vector field  $m : \mathbb{R}^d \rightarrow \mathbb{R}^d$  specified by  
 184 discretizing (2) and (3). Under deterministic dynamics,  $\mathbf{X}(t) \rightarrow \mathbf{a}$  or  $\mathbf{X}(t) \rightarrow \mathbf{b}$  as  $t \rightarrow \infty$  depending  
 185 on initial conditions. The addition of white noise changes the system into an Itô diffusion

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

186 where  $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$  imparts a correlation structure to the vector  $\mathbf{W}(t) \in \mathbb{R}^d$  of independent  
 187 standard white noise processes. We design  $\sigma$  to be a low-rank, constant matrix that adds spatially  
 188 smooth stirring to only the zonal wind  $U$  (not the streamfunction  $\Psi$ ) and respects boundary  
 189 conditions at the bottom and top of the stratosphere. We simulate the model using the Euler-  
 190 Maruyama method: in a timesetep  $\delta t$ , after a deterministic forward Euler step we add the stochastic  
 191 perturbation to zonal wind on large vertical scales

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

192 where  $\eta_k$  ( $k = 0, 1, 2$ ) are independent unit normal samples. We set the magnitude of  $\sigma$  by

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

193  $\sigma_U$  technically has units of  $(L/T)/T^{1/2}$ , where the square-root of time comes from the quadratic  
 194 variation of the Wiener process. It is best interpreted in terms of the daily root-mean-square  
 195 velocity perturbation of 1.0 m/s.

196 A long stochastic simulation of the model reveals metastability, with the system tending to  
 197 remain close to one fixed point for a long time before switching quickly to the other, as shown  
 198 by the timeseries of  $U(30\text{ km})$  in panel (c) of Figure 1. We display the zonal wind  $U$  at 30 km  
 199 following Christiansen (2000), because this is about where zonal wind strength is minimized in  
 200 the weak vortex. While the two regimes are clearly associated with the two fixed points, they are  
 201 better characterized by extended *regions* of state space with strong and weak vortices. We thus  
 202 define the two metastable subsets of  $\mathbb{R}^d$

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

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

203 This straightforward definition roughly follows the convention of Charlton and Polvani (2007),  
 204 which defines an SSW as a reversal of zonal winds at 10 hPa. In the Holton-Mass model, where  
 205  $z = 0$  at the tropopause, this translates to  $z = -7\text{ km} \ln(10/1000) - 10\text{ km} = 22.2\text{ km}$ , but we have  
 206 adjusted the specific altitude here to 30 km where the zonal wind reduction is most drastic. There is  
 207 lively debate around the definition of SSW events (e.g., Butler et al. 2015), with different thresholds  
 208 leading to different statistics. The details are affected by the definition in our analysis, but the results  
 209 are qualitatively similar over a wide range. Our method is equally applicable to any definition,  
 210 and so to illustrate we choose one that enjoys broad acceptance. Incidentally, the analysis tools we  
 211 present may be helpful in distinguishing predictability properties between different definitions.

212 The green highlights in Figure 1 (c) begin precisely when the system exits the  $A$  region bound  
 213 for  $B$ , and end when the system enters  $B$ . The orange highlights start when the system leaves  $B$   
 214 bound for  $A$ , and end when  $A$  is reached. Note that  $A \rightarrow B$  transitions are much shorter in duration  
 215 than  $B \rightarrow A$  transitions. Figure 1 (d) shows the same paths, but viewed in the space  $(|\Psi|, U)$  at 30  
 216 km. The  $A \rightarrow B$  and  $B \rightarrow A$  transitions are again highlighted in green and orange respectively,

217 showing geometrical differences between the two directions. We will refer to the  $A \rightarrow B$  transition  
 218 as an SSW event, even though it is more accurately a transition between climatologies according  
 219 to the Holton-Mass interpretation. The  $B \rightarrow A$  transition is a vortex restoration event. Our focus  
 220 in this paper is on predicting transition events and monitoring their progress in a principled way.  
 221 In the next section we explain the formalism for doing so.

### 222 3. Theory and computation

#### 223 a. Definitions

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

$$\tau_{D^c} = \min\{t \geq 0 : \mathbf{X}(t) \in D^c\} \quad (9)$$

230 Because of stochastic forcing (which in practice arises from unobserved variables), the first-exit  
 231 time is a random variable, formally called a “stopping time” (Oksendal 2003; Durrett 2013),  
 232 meaning measurable from the history of  $\mathbf{X}(t)$ . The first-exit location  $\mathbf{X}(\tau_{D^c})$  is itself a random  
 233 variable which importantly determines how the system exits  $D$ : either  $\mathbf{X}(\tau_{D^c}) \in A$ , meaning the  
 234 vortex restores to radiative equilibrium, or  $\mathbf{X}(\tau_{D^c}) \in B$ , meaning the vortex breaks down into  
 235 vacillation cycles. A fundamental goal of probabilistic forecasting is to determine the probabilities

236 of these two events, which naturally leads to the definition of the (forward) committor function

$$q^+(\mathbf{x}) = \begin{cases} \mathbb{P}_{\mathbf{x}}\{\mathbf{X}(\tau_{D^c}) \in B\} & \mathbf{x} \in D = (A \cup B)^c \\ 0 & \mathbf{x} \in A \\ 1 & \mathbf{x} \in B \end{cases} \quad (10)$$

237 where the subscript  $\mathbf{x}$  indicates that the probability is conditional on a fixed initial condition  
 238  $\mathbf{X}(0) = \mathbf{x}$ , i.e.,  $\mathbb{P}_{\mathbf{x}}\{\cdot\} = \mathbb{P}\{\cdot | \mathbf{X}(0) = \mathbf{x}\}$ . (The superscript “+” distinguishes the forward committor  
 239 from the *backward committor*, an analogous quantity for the time-reversed process which we do  
 240 not use in this paper.) Throughout, we will use capital  $\mathbf{X}(t)$  to denote a stochastic process, and  
 241 lower-case  $\mathbf{x}$  to represent a specific point in state space, typically an initial condition, i.e.,  $\mathbf{X}(0) = \mathbf{x}$ .  
 242 Both are  $d = 75$ -dimensional vectors. The boundary conditions in Equation (10) naturally extend  
 243 the definition of  $q^+$  in  $D$ : if the vortex starts out very weak,  $\mathbf{x}$  is close to set  $B$  and the system will  
 244 probably land in  $B$  next, making  $q^+(\mathbf{x}) \approx 1$ . If it starts out strong and close to  $A$ , it will most likely  
 245 restore to  $A$  next, making  $q^+(\mathbf{x}) \approx 0$ . The committor is clearly a function of initial condition  $\mathbf{x}$ , but  
 246 assuming the process is Markovian, it does not depend on the history of the system that led it to  $\mathbf{x}$   
 247 in the first place.

248 Another important forecasting quantity is the lead time to the event of interest. While the forward  
 249 committor reveals the probability of experiencing vortex breakdown *before* returning to a strong  
 250 vortex, it does not say how long either event will take in absolute terms. Furthermore, even if  
 251 the vortex is restored first, how long will it be until the next SSW does occur? The time until the  
 252 next SSW event is denoted  $\tau_B$ , again a random variable, whose distribution depends on the initial  
 253 condition  $\mathbf{x}$ . We call  $\mathbb{E}_{\mathbf{x}}[\tau_B]$  the *mean first passage time* (MFPT) to  $B$ . Conversely, we may ask  
 254 how long a vortex disturbance will persist before normal conditions return; the answer (on average)  
 255 is  $\mathbb{E}_{\mathbf{x}}[\tau_A]$ , the mean first passage time to  $A$ . Dissecting the expectations further, we may condition

256  $\tau_B$  on the event that an SSW is coming before the strong vortex returns, leading to the conditional  
257 first passage time  $\mathbb{E}_{\mathbf{x}}[\tau_B | \tau_B < \tau_A]$ , which in some sense quantifies the suddenness of SSW.

258 All of these quantities can, in principle, be estimated by collecting averages over very long  
259 simulations. For example, to estimate the committor at a given  $\mathbf{x}$ , one can shoot  $N$  trajectories  
260 starting from  $\mathbf{x}$  and count the numbers  $N_A$  and  $N_B$  hitting  $A$  and  $B$  first. Then  $N_A/N$  will be an  
261 estimate for the committor at  $\mathbf{x}$ . The mean first passage time can be estimated using these same  
262 sampled trajectories. But this direct method can be prohibitively expensive, especially if applied  
263 to many points all over state space. By definition, transitions between  $A$  and  $B$  are infrequent.  
264 Therefore, if starting from  $\mathbf{x}$  far from  $B$ , then a huge number of sampled trajectories ( $N$ ) will be  
265 required to observe even a small number ending in  $B$  ( $N_B$ ). Likewise, transition path statistics such  
266 as return times can, in principle, be computed from an extremely long model run, but in most cases  
267 of interest this direct simulation approach will not be feasible.

268 In Subsection 3(b) we will write these forecasts in a single general form, and describe a com-  
269 putational approach to compute them using only a data set of short forward model integrations.  
270 The method is called the Dynamical Galerkin Approximation (DGA), introduced in Thiede et al.  
271 (2019). It takes advantage of the Feynman-Kac formula (Oksendal 2003), recasting conditional  
272 expectations as PDE problems over state space. These equations are *local* and thus approximable  
273 by short trajectories. We perform these calculations on the Holton-Mass model and present the  
274 results in Section 4. In Section 5(a) we describe a statistical analysis to aid interpretation of the  
275 estimated committor.

### 276 *b. Dynamical Galerkin Approximation*

277 In this section we describe the methodology, which involves some technical results from stochastic  
278 processes and measure theory. The forecast functions described above—committors and MFPTs—

279 can all be written as conditional expectations of the form

$$F(\mathbf{x}) = \mathbb{E}_{\mathbf{x}} \left[ G(\mathbf{X}(\tau)) e^{-\int_0^\tau V(\mathbf{X}(s)) ds} + \int_0^\tau H(\mathbf{X}(s)) e^{-\int_0^s V(\mathbf{X}(r)) dr} ds \right] \quad (11)$$

280 where again the subscript  $\mathbf{x}$  denotes conditioning on  $\mathbf{X}(0) = \mathbf{x}$ ;  $G, H$  and  $V$  are arbitrary known  
 281 functions over  $\mathbb{R}^d$ ; and  $\tau$  is a stopping time, specifically a first-exit time like Equation (9) but  
 282 possibly with  $D$  replaced by another set. To see that the forward committor takes on this form,  
 283 set  $G(\mathbf{x}) = \mathbb{1}_B(\mathbf{x})$  (one on set  $B$  and zero everywhere else),  $V = 0$ ,  $H = 0$ , and  $\tau = \tau_{D^c}$ . Then  
 284  $F(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[\mathbb{1}_B(\mathbf{X}(\tau))] = \mathbb{P}_{\mathbf{x}}\{\mathbf{X}(\tau_{D^c}) \in B\} = q^+(\mathbf{x})$ . For the mean first passage time to  $B$ , set  $\tau = \tau_B$ ,  
 285  $G = 0$ ,  $V = 0$ , and  $H = 1$ . Then  $F(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[\int_0^{\tau_B} dt] = \mathbb{E}_{\mathbf{x}}[\tau_B]$ . For the conditional first passage time  
 286 to  $B$ , set  $G = \mathbb{1}_B$ ,  $\tau = \tau_{D^c}$ ,  $V = -\lambda$  (a constant) and  $H = 0$ . Then the expectation can be computed  
 287 in two steps, by computing and differentiating a moment-generating function:

$$\begin{aligned} F(\mathbf{x}; \lambda) &= \mathbb{E}_{\mathbf{x}}[\mathbb{1}_B(\mathbf{X}(\tau_{D^c})) e^{\lambda \tau_{D^c}}] \\ \frac{1}{q^+(\mathbf{x})} \frac{\partial}{\partial \lambda} F(\mathbf{x}; 0) &= \frac{\mathbb{E}_{\mathbf{x}}[\tau_{D^c} \mathbb{1}_B(\mathbf{X}(\tau_{D^c}))]}{\mathbb{E}_{\mathbf{x}}[\mathbb{1}_B(\mathbf{X}(\tau_{D^c}))]} \\ &= \mathbb{E}_{\mathbf{x}}[\tau_B | \tau_B < \tau_A] \end{aligned} \quad (12)$$

288 Note that the event  $\tau_B < \tau_A$  is equivalent to  $\mathbf{X}(\tau_{D^c}) \in B$ , where again  $D^c = A \cup B$ . We approximate  
 289  $\partial F / \partial \lambda$  with a centered finite difference, after computing  $F(\mathbf{x}; \lambda)$  for several  $\lambda$  in the neighborhood  
 290 of zero. In principle we could compute higher moments in the same way and get a more detailed  
 291 understanding of the conditional passage time distribution. Alternatively we could estimate  $F(\mathbf{x}; \lambda)$   
 292 for a large range of  $\lambda$  and recover the distribution with a Laplace transform.

293 More generally, the function  $G$  is chosen by the user to quantify risk at the terminal time  $\tau$ ; in  
 294 the case of the forward committor, that risk is binary, with an SSW representing a positive risk and  
 295 a radiative vortex no risk at all. The function  $H$  is chosen to quantify the risk accumulated up until

296 time  $\tau$ , which might be simply an event's duration, but other integrated risks may be of more interest  
 297 for the application. For example, one could express the total thermal energy absorbed by the polar  
 298 vortex by setting  $H = \overline{v'T'}$ , or the momentum lost by the vortex by setting  $H(\mathbf{x}) = U(\mathbf{a}) - U(\mathbf{x})$ , or  
 299 a vertically integrated version. Using the moment generating function in (12), one can compute  
 300 not only means but higher moments of such integrals by expressing the risk with  $V$ .

301 Let us now describe how to numerically compute  $F(\mathbf{x})$  of the form (11) with short trajectories,  
 302 starting with the special case of the forward committor and then generalizing. Consider starting a  
 303 random trajectory at  $\mathbf{x} = \mathbf{X}(0) \in D = (A \cup B)^c$  and evolving it for a short time  $\Delta t$ . Its probability  
 304 of reaching  $B$  first,  $q^+(\mathbf{x})$ , is simply the probability that it reaches  $B$  first starting from  $q^+(\mathbf{X}(\Delta t))$   
 305 instead, averaged over all possible  $\mathbf{X}(\Delta t)$  (ignoring momentarily the small probability that  $A \cup B$  is  
 306 reached before time  $\Delta t$ ). That is,  $q^+(\mathbf{x}) \approx \mathbb{E}_{\mathbf{x}}[q^+(\mathbf{X}(\Delta t))] =: \mathcal{T}^{\Delta t} q^+(\mathbf{x})$ . The operator  $\mathcal{T}^{\Delta t}$  is known  
 307 as the (stochastic) transition operator, which maps a function on state space to the expectation of  
 308 that function at a future time. We could furthermore divide by  $\Delta t$  and take the limit  $\Delta t \rightarrow 0$ ,  
 309 eliminating the event  $\tau_{D^c} < \Delta t$ , and obtain the Kolmogorov Backward PDE (e.g., Oksendal 2003;  
 310 Weinan et al. 2019). Instead, to represent our numerical method more directly, we implement  
 311 a purely finite-time approach from Strahan et al. (2020): artificially halt the dynamics upon first  
 312 arrival in  $D^c = A \cup B$  and modify the equation to

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}}[q^+(\mathbf{X}(\Delta t \wedge \tau_{D^c}))] - q^+(\mathbf{x}) \\ & = (\mathcal{T}_{D^c}^{\Delta t} - 1)q^+(\mathbf{x}) = 0 \end{aligned} \tag{13}$$

313 where  $\Delta t \wedge \tau_{D^c} := \min(\Delta t, \tau_{D^c})$  and  $\mathcal{T}_{D^c}^{\Delta t}$  is a "stopped" transition operator. This equation holds  
 314 for  $\mathbf{x} \in D$ , and comes with the boundary condition  $q^+(\mathbf{x}) = \mathbb{1}_B(\mathbf{x})$  for  $\mathbf{x} \in A \cup B$ . Applying similar  
 315 logic to the mean first passage time to  $B$ , let  $\mathbf{x} \in B^c$ , denote  $m_B(\mathbf{x}) = \mathbb{E}_{\mathbf{x}}[\tau_B]$  and observe that  $\tau_{D^c}$

316 decreases by  $\Delta t \wedge \tau_{B^c}$  during the short timespan. So for all  $\mathbf{x} \in B^c$ ,

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}}[m_B(\mathbf{X}(\Delta t \wedge \tau_B^+))] - m_B(\mathbf{x}) \\ &= (\mathcal{T}_B^{\Delta t} - 1)m_B(\mathbf{x}) = -\mathbb{E}_{\mathbf{x}}[\Delta t \wedge \tau_B] \end{aligned} \quad (14)$$

317 with the boundary condition  $m_B(\mathbf{x}) = 0$  for  $\mathbf{x} \in B$ . Now in the general case, let  $D$  stand in for the  
 318 relevant region of state space and let  $G, H$  and  $V$  be arbitrary. The corresponding operator equation  
 319 is

$$\begin{aligned} & (\mathcal{T}_{D^c}^{\Delta t} - 1)F(\mathbf{x}) - \mathbb{E}_{\mathbf{x}} \left[ \int_0^{\Delta t \wedge \tau_{D^c}} V(\mathbf{X}(t))F(\mathbf{X}(t)) dt \right] \\ &= -\mathbb{E}_{\mathbf{x}} \left[ \int_0^{\Delta t \wedge \tau_{D^c}} H(\mathbf{X}(t)) dt \right] \end{aligned} \quad (15)$$

320 for  $\mathbf{x} \in D$ , with boundary condition  $F(\mathbf{x}) = G(\mathbf{x})$  for  $\mathbf{x} \in D^c$ . This linear equation comes from  
 321 Dynkin's formula, an integrated version of the Feynman-Kac; see Oksendal (2003); Karatzas and  
 322 Shreve (1998); Weinan et al. (2019) theoretical background. The remarkable aspect of this formula  
 323 is that while  $F$  is an expectation over paths going all the way to the boundary  $D^c$  (a strong or  
 324 weak vortex), it obeys a *local* equation with expectations over short trajectories of length  $\Delta t$ . By  
 325 collecting many short-trajectory samples, we can compute statistical properties of the event without  
 326 ever actually observing one happen in simulation. Note that (15) reduces to (13) with  $V = H = 0$   
 327 and (14) with  $V = 0, H = 1$ .

328 Like a PDE with a high dimensional independent variable space, Equation (15) cannot be solved  
 329 using any classical discretization of the possible values of  $\mathbf{x}$ . Successful approaches will involve  
 330 a representation of the solution,  $F$ , suitable for the high dimensional setting, i.e. representations  
 331 of the type commonly employed for machine learning tasks. The DGA method, in particular,  
 332 consists of expanding the unknown function  $F$  in a "data-informed" basis (to be specified later).  
 333 The expectations in Equation (15) are estimated by launching short trajectories from all over state

space. Finally, a finite system of equations is solved for the unknown coefficients in the basis expansion of  $F$ , in effect stitching together information from all trajectories at once.

We can express the essential idea using the example of Equation (13) for  $q^+(\mathbf{x})$ , while the supplement contains a more general version. We first homogenize the boundary conditions with a guess function  $\hat{q}^+(x)$  that obeys the boundary conditions  $\hat{q}^+|_A = 0$ ,  $\hat{q}^+|_B = 1$ , and let  $r(x) = q^+(x) - \hat{q}^+(x)$ , so that  $r$  obeys homogeneous Dirichlet conditions and satisfies

$$(\mathcal{T}_{D^c}^{\Delta t} - 1)r(x) = -(\mathcal{T}_{D^c}^{\Delta t} - 1)\hat{q}^+(x) \quad (16)$$

We next expand  $r$  in a finite-dimensional basis of functions  $\{\phi_1, \dots, \phi_M\}$  with unknown coefficients  $c_j$ :  $r(x) = \sum_{j=1}^M c_j(r)\phi_j(x)$ . Each  $\phi_j$  obeys the homogeneous boundary conditions. Finally, we take the inner product of both sides with  $\phi_i$ , with respect to some measure  $\mu$ , to produce a system of  $M$  linear equations

$$\sum_{j=1}^M \langle \phi_i, (\mathcal{T}_{D^c}^{\Delta t} - 1)\phi_j \rangle_{\mu} c_j(r) = -\langle \phi_i, (\mathcal{T}_{D^c}^{\Delta t} - 1)\hat{q}^+ \rangle_{\mu} \quad i = 1, \dots, M \quad (17)$$

These inner products are intractable integrals over high-dimensional state space, but can be approximated using Monte Carlo integration. If  $\mathbf{X}$  is an  $\mathbb{R}^d$ -valued random variable distributed according to  $\mu$ , and we have access to random samples  $\{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ , the law of large numbers gives

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N f(\mathbf{X}_n) = \int_{\mathbb{R}^d} f(\mathbf{x}) \mu(d\mathbf{x}) \quad (18)$$

This is where the short trajectory data enters the picture. We generate a dataset of length- $\Delta t$  trajectories  $\{\mathbf{X}_n(t) : 0 \leq t \leq \Delta t, n = 1, \dots, N\}$ . These short trajectories might enter  $A$  or  $B$  before time  $\Delta t$ , and to account for this we also store the stopping times  $\Delta t \wedge \tau_{n,D^c}$ . The  $\mathbf{X}_n(0)$ 's are sampled from an arbitrary measure  $\mu$ , called the sampling measure, which is determined by the sampling procedure for initial points. For example, if points are selected randomly from a long

trajectory,  $\mu \approx \pi$  (the steady-state probability density) by ergodicity. However, we may choose  $\mu$  so that many samples appear in regions of particular interest, such as transition regions far away from  $A$  and  $B$  and to which  $\pi$  assigns very little probability. Once the dataset is generated, we use  $\mu$  as the reference measure for the inner products in (17), allowing us to approximate them with Monte Carlo integration. For example,

$$\begin{aligned} & \langle \phi_i, (\mathcal{T}_{D^c}^{\Delta t} - 1)\phi_j \rangle_{\mu} \\ & \approx \frac{1}{N} \sum_{n=1}^N \phi_i(\mathbf{X}_n(0)) [\phi_j(\mathbf{X}_n(\Delta t \wedge \tau_{n,D^c})) - \phi_j(\mathbf{X}_n(0))] \end{aligned} \quad (19)$$

We can similarly estimate any expectation of the form (15) using different basis functions adapted to the specific region of interest.

The formulation above works for any class of basis functions that becomes increasingly expressive as the library grows, capable of estimating any function of interest. However, with a finite truncation, choosing the basis functions is a crucial ingredient of DGA, greatly impacting the efficiency and accuracy of the results. In our current study, we restrict to the 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 disjoint partition of state space. In practice we will construct these sets by clustering data. This basis set construction is borrowed from common practice in the computational statistical mechanics community for building a Markov State Model (MSM) (Frank and Fischer 2008; Pande et al. 2010; Bowman et al. 2013; Chodera and Noé 2014). MSMs are a dimensional reduction technique that has also been used in conjunction with analysis of metastable transitions, primarily in protein folding dynamics (Noé et al. 2009) and were recently used to study ocean circulation in Miron et al. (2021). DGA can be viewed as an extension of MSMs, though, rather than producing any reduced complexity model, the explicit goal in DGA is the estimation of specific functions as in Equation (11). The supplement

372 spells out DGA in considerably greater detail, which may be more helpful to view after seeing the  
373 forthcoming results.

## 374 4. Methods

375 In this section we explain our specific application of the Dynamical Galerkin Approximation  
376 (DGA) method (Thiede et al. 2019; Strahan et al. 2020) to the Holton-Mass model and validate  
377 our results empirically from simulation.

### 378 *a. Data generation*

379 There are many possible ways to choose starting points for the short trajectories. Whatever  
380 procedure we use will induce a *sampling measure*  $\mu$  on state space.  $\mu(\mathbf{x})$  is a probability density  
381 that specifies the expected number of starting points per unit volume in the region of state space near  
382  $\mathbf{x}$ . This is a natural reference measure for the Monte Carlo inner products described in Subsection  
383 3(b). Because  $\mu$  has minimal requirements, the user is afforded great flexibility in sampling the  
384 data. How to efficiently generate maximally informative data is an active and nontrivial research  
385 question, but a few heuristics are obvious. In a metastable system, setting  $\mu = \pi$  would be a poor  
386 choice, because the data would be strongly concentrated in the immediate neighborhoods of  $A$  and  
387  $B$ , whereas the regions of primary interest are the transition regions somewhere in between  $A$  and  $B$ .  
388 Different physical observables, such as the Eliassen-Palm flux, may be important prior candidates  
389 for their predictive power, and we might like to seed data samples uniformly over a certain range  
390 of that variable. On the other hand, the samples should fall within a physically realistic region of  
391 state space, not just any point in  $\mathbb{R}^{75}$ . To see why, recall that the last 25 entries of the state vector  
392  $\mathbf{X}$  represent the velocity field at discretized vertical levels from  $z = 0 \text{ km}$  to  $z = 70 \text{ km}$ . Because

393 velocity is a continuous function of altitude, adjacent entries should be close together, which is not  
394 at all guaranteed for a randomly chosen 75-dimensional vector.

395 Because a goal of this article is to demonstrate interpretable results of rare event analysis on  
396 a climate model, we choose an easy, probably suboptimal sampling strategy. We defer opti-  
397 mization to later work, perhaps for a more expensive model that demands it. We define our  
398 sampling distribution as the equilibrium distribution, re-weighted to be uniform over the space  
399  $(U(30\text{ km}), |\Psi|(30\text{ km}))$  within the bounds realized by the control simulation, which are approx-  
400 imately  $-30\text{ m/s} \leq U(30\text{ km}) \leq 70\text{ m/s}$  and  $0\text{ m}^2/\text{s} \leq |\Psi|(30\text{ km}) \leq 2 \times 10^7\text{ m}^2/\text{s}$ . Without direct  
401 access to the equilibrium distribution, we approximate it by running a very long trajectory of  
402 500,000 days, producing many transitions like those shown in Figure 1, with an Euler-Maruyama  
403 timestep of 0.005 days (for comparison, a single transition event takes on the order of 100 days).  
404 We acknowledge this is cheating on our claim to only use short trajectories; however, we use  
405 the long simulation only to seed the initial conditions for those short trajectories, as well as to  
406 empirically validate the results of DGA later on. This way we can emphasize the power of DGA  
407 itself, which will motivate more efficient upstream data generation methods. Alternatives exist for  
408 sampling state space thoroughly without a long simulation, for example trajectory-splitting (e.g.,  
409 L'Ecuyer et al. 2007). One could initialize trajectories in one of the metastable sets, say in  $A$  di-  
410 rectly on the fixed point  $\mathbf{a}$ , and integrate the trajectories for a short time to explore the surrounding  
411 region. These new data points can be used as initial conditions for the next round of simulation, at  
412 each stage exploring a wider region of state space until the bulk of the attractor is covered. This  
413 initialization procedure may require a long total simulation time, but is parallelizable. We will  
414 explore and optimize such methods in future work with more sophisticated models, where efficient  
415 initialization is more critical. For now we settle for initial data points from a long simulation.

416 After downsampling the long simulation to a resolution of 0.5 days, we sample snapshots from  
 417 the trajectory, reweighted to induce a uniform distribution on the space  $(U(30\text{ km}), |\Psi(30\text{ km})|)$ .  
 418 Specifically, we compute a discrete histogram over the two-dimensional space and weight each  
 419 sample by the inverse of its density on that histogram. We collect  $N = 1 \times 10^6$  snapshots  $\{\mathbf{X}_n(0)\}$   
 420 directly from the long simulation, and then launch independent (hence completely parallelizable)  
 421 short 10-day trajectories from each, to obtain the short trajectory database  $\{\mathbf{X}_n(t) : 0 \leq t \leq \Delta t, n =$   
 422  $1 \dots, N\}$ . Afterward we identify the first-entry times to  $D^c$  for each trajectory, called  $\tau_{n,D^c}$ . This  
 423 strategy is straightforward and guarantees that  $\mu$  gives substantial probability to candidate transition  
 424 regions and that only physically reasonable points are sampled.

425 *b. Computation and validation*

426 The partition  $\{S_1, \dots, S_M\}$  to build the basis function library  $\{\mathbb{1}_{S_j}(\mathbf{x})\}_{n=1}^N$  should be chosen with  
 427 a number of considerations in mind. The partition elements should be small enough to accurately  
 428 represent the functions they are used to approximate, but large enough to contain sufficient data to  
 429 robustly estimate transition probabilities. We form these sets by a hierarchical modification of  $K$ -  
 430 means clustering on  $\{\mathbf{X}_n(0)\}_{n=1}^N$ .  $K$ -means is a robust method that can incorporate new samples by  
 431 simply identifying the closest centroid, and is commonly used in molecular dynamics (Pande et al.  
 432 2010). However, straightforward application of  $K$ -means, as implemented in the `scikit-learn`  
 433 software (Pedregosa et al. 2011), can produce a very imbalanced cluster size distribution, even with  
 434 empty clusters. This leads to unwanted singularities in the constructed Markov matrix. To avoid  
 435 this problem we cluster hierarchically, starting with a coarse clustering of all points and iteratively  
 436 refining the larger clusters, at every stage enforcing a minimum cluster size, until we have the  
 437 desired number of clusters ( $M$ ). After clustering on the initial points  $\{\mathbf{X}_n(0)\}$ , the other points  
 438  $\{\mathbf{X}_n(t), 0 < t \leq \Delta t\}$  are placed into clusters using an address tree produced by the  $K$ -means cluster

439 hierarchy. To guarantee that  $D$  and  $D^c$  consist exactly of a union of subsets, we cluster points in  
440  $D$  and  $D^c$  separately, with a number of clusters proportional to the number of points therein. (We  
441 remind the reader that the domain and boundary depend on which quantity of interest is being  
442 computed. For the forward and backward committor,  $D^c$  consists of  $A$  and  $B$ , which are defined  
443 *a priori* by thresholds of  $U(30\text{ km})$ .) The total number of clusters is fixed to  $M = 1500$ . When  
444 doing out-of-sample extension on a point  $z$ , we first identify whether  $z \in D$  or  $D^c$ , and assign it to  
445 a cluster accordingly.

446 Figure 2 demonstrates the accuracy of the calculated forward committor and mean first passage  
447 time to  $B$  by taking advantage of the long trajectory from which we sampled the short trajectories.  
448 We divide the interval  $(0, 1)$  into 20 bins, and identify for each interval  $(\zeta_1, \zeta_2)$  which data points  
449  $\{\mathbf{X}_n(0) : \zeta_1 < q^+(\mathbf{X}_n(0)) < \zeta_2\}$  were en route to the vacillating regime at the instant they were  
450 selected from the long simulation. If the committor is computed accurately, the proportion of  
451 data points headed to  $B$  should fall in the interval  $(\zeta_1, \zeta_2)$ . For example, about 20-25% of data  
452 points  $\mathbf{X}_n(0)$  whose committor is calculated to be within  $(0.2, 0.25)$  should be headed to set  
453  $B$ . Analogously, we expect rain 20% of the time the National Weather Service forecasts a 20%  
454 chance of rain. This is a very coarse measure of accuracy, and only a necessary condition, but the  
455 strong empirical match shown in the scatter plots of Figure 2 gives us confidence in our numerical  
456 results. The mean first passage time calculation is evaluated similarly: for all data points  $\mathbf{X}_n(0)$   
457 with the estimated  $m_B(\mathbf{X}_n(0))$  in a certain range  $(t_1, t_2)$ , we average the true first-passage time  
458 observed from the long trajectory. The match is quite good up until very long lag times, where  
459 DGA underestimates the long tail. The accuracy of committors and passage times improve as the  
460 dataset grows and clusters are refined. More sophisticated basis sets and sampling methods may  
461 significantly improve the convergence rate.

462 The committor and first passage time relate to the weather forecasting problem of predicting the  
 463 next rare event given the current initial condition. However, they can also characterize the polar  
 464 vortex climatology, meaning its average behavior over very long time periods as pertains to  $A$  and  
 465  $B$ . To wit, how much does the system “prefer” to be in a weak or strong state, as measured by the  
 466 fraction of time it spends in either? This can be quantified by the steady state distribution (also called  
 467 the invariant or stationary measure)  $\pi(\mathbf{x})$ , the probability distribution function produced by binning  
 468 data points from a very long simulation. Figure 3 illustrates that the metastable Holton-Mass model  
 469 has a starkly bimodal distribution, with the system tending to spend a long time in state  $A$  or  $B$   
 470 before occasionally switching quickly to the other state. We have estimated  $\pi$  here using a variation  
 471 on the DGA recipe described above. The details of the calculation can be found in the supplement.  
 472 We have projected  $\pi$  onto the two-dimensional subspace  $(|\Psi(30\text{ km})|, U(30\text{ km}))$  on a log scale,  
 473 along with a one-dimensional projection onto the latter coordinate  $U(30\text{ km})$  on a linear scale. The  
 474 preferences for  $A$  and  $B$  can be quantitatively compared by the fraction of time spent inside each set,  
 475 as well as the fraction of time spent between the two sets but destined for either one. These ergodic  
 476 averages can be found by averaging the forward committor over different regions of state space.  
 477 For example, the fraction of time spent inside  $A$  is  $\int_A \pi(d\mathbf{x}) = \int_{\mathbb{R}^d} \mathbb{1}_A(\mathbf{x})\pi(d\mathbf{x}) = \langle \mathbb{1}_A \rangle_\pi$ . Similarly,  
 478 the fraction of time spent inside  $B$  is  $\langle \mathbb{1}_B \rangle_\pi$ ; the fraction spent outside  $A$  and  $B$  but destined for  $B$   
 479 is  $\langle \mathbb{1}_{(A \cup B)^c} q^+ \rangle_\pi$ ; and the fraction spent outside  $A$  and  $B$  but destined for  $A$  is  $\langle \mathbb{1}_{(A \cup B)^c} (1 - q^+) \rangle_\pi$ .  
 480 Table 1 displays these fractions calculated from DGA and empirically from the long trajectory.  
 481 The time spent either in  $A$  or destined for  $A$  (the first two rows) is about equal to the time spent in  
 482 or destined for  $B$  (last two rows). However, more time is spent destined for  $B$  than strictly inside  
 483  $B$ , as vacillation cycles often increase the zonal wind above 1.75 m/s before it dips back down.  
 484 Furthermore, Figure 3 shows a higher and narrower peak in the  $A$  regime. We interpret that a

485 strong vortex is much less variable than a weak vortex, which is consistent with the vacillation  
486 cycles that characterize the latter.

## 487 **5. Results and Discussion**

488 Our analysis can be roughly divided into two parts. First, from a forecasting perspective, we  
489 demonstrate that the committor is more robust than naïve proxies from the model as a leading  
490 indicator of an oncoming SSW. We also find a low-rank representation of the committor in terms  
491 of the system’s basic observables using a sparsity-promoting LASSO regression (Tibshirani 1996).  
492 Second, we quantitatively relate the *risk* of an oncoming event with the *lead time* to the event, an  
493 important consideration in extreme weather prediction.

### 494 *a. The committor as an early warning*

495 Operational forecasting requires continuous updating of probabilities from incoming observa-  
496 tions, which provide only partial information on the state of the atmosphere. The choice of which  
497 observables to monitor is constrained by measurement capabilities, but is also informed by pre-  
498 diction efficacy; we desire warning signs that are highly correlated with the event and occur as  
499 early as possible to give some buffer time to brace for impacts. Figure 4 visually demonstrates the  
500 advantage of considering the committor as a forecasting metric compared to two other observables:  
501 zonal wind  $U$  and meridional eddy heat flux  $\overline{v'T'}$ , both measured at the same altitude of 30 km.  
502 We have extracted a typical complete SSW event ( $A \rightarrow B$  transition path) from the long simulation  
503 and plotted a timeseries of the observables on a common time axis. The time  $t = 0$  corresponds  
504 to the central date of a warming event, the moment when the system first enters set  $B$ , with zonal  
505 wind at 30 km dropping below the threshold of 1.75 m/s. The committor timeseries (Figure 4a) is  
506 estimated by nearest neighbor interpolation from the dataset  $\{\mathbf{X}_n(0)\}_{n=1}^N$ .

507 The committor curve timeseries first exceeds the threshold of 0.5 around 27 days before the  
508 event while rising sharply in a roughly S-shaped curve. A perfect committor-measuring instrument  
509 would be sending a strong signal of increasing risk at that time. Compare this with  $U(t)$ , which is  
510 plateauing, or very gradually decreasing, around  $40\text{ m/s}$  when the threshold  $q^+ = 0.5$  is crossed. The  
511 apparently mild behavior belies the rapid increase in SSW risk shown by the committor timeseries.  
512 The dramatic drop in zonal wind occurs well after the committor exceeds 0.5, and so a reading  
513 of  $U$  directly would not give a strong warning sign until late in the progress of transition. One  
514 could write the committor as an approximate function of  $U(30\text{ km})$ , which is plotted in Figure 5(a)  
515 as explained below, and would find that the  $U$ -level corresponding to  $q^+ = 0.5$  is around  $37\text{ m/s}$ .  
516 Unfortunately,  $U(30\text{ km})$  does not drop below  $37\text{ m/s}$  until the SSW is 12 days away, providing  
517 much less lead time than if the full committor were known. The considerable gap in prediction  
518 date is shown by a blue strip. Meanwhile, the heat flux over time plotted in panel (c) suffers the  
519 same deficiency as a predictor, having an analogous threshold of  $1.2 \times 10^{-6}\text{ K} \cdot \text{m/s}$ . The  $\overline{v'T'}$  level  
520 hardly budes while critical preconditions are falling into place, and only after the die is already  
521 cast in favor of a SSW does the heat flux rise sharply. A monitoring system based on heat flux  
522 alone would be very ill-informed about the risk of impending SSW event. While heat flux is a  
523 dynamically consequential quantity for describing the evolution of an SSW, this does not directly  
524 translate into good predictive properties. These prediction gaps are typical: over many simulated  
525 transitions, the average delay between  $q^+(\mathbf{X}(t))$  clearing 0.5 for the last time the other observables  
526 clearing their thresholds for the last time are 9.1 days for  $U(30\text{ km})$  and 9.8 days for  $\overline{v'T'}(30\text{ km})$ .

527 We use the caveat “directly” because the possibility remains that the vertical scale in Figures 4  
528 (b-c) unfairly downplay the predictive power of zonal wind and heat flux. Perhaps they could be  
529 very robust predictors, if examined on the right scale and with appropriate (possibly nonlinear)  
530 transformations. Calculating such a transformation on theoretical grounds alone would be a

531 daunting task, especially in light of stochastic perturbations. But even if this were possible, at best  
532 this calculation would approximate nothing other than the forward committor itself. Furthermore,  
533 we incorporate all state variables at once into the committor calculation, which is at least as flexible  
534 as considering heat flux or zonal wind alone. Nonetheless, for the sake of dynamical transparency  
535 and practical observational constraints, it would be helpful to have a parsimonious representation  
536 of the committor in terms of a small number of state variables, if possible. We pursue this prospect  
537 in the following subsection.

538 *b. Sparse representation of the committor*

539 The committor's superiority as a probabilistic forecast is not surprising, because it is built into  
540 the definition. The committor combines information from every degree of freedom in just the right  
541 way to give the probability of next hitting  $B$  rather than  $A$ . However, these degrees of freedom  
542 may not all be "observable" in a practical sense, given the sparsity and resolution limits of weather  
543 sensors. It is therefore important to ask: what is the best possible estimate of the committor given  
544 an observed subset of state variables? A related question arises in the design of observational  
545 systems: which variables should be measured to optimally estimate the committor, under cost and  
546 engineering constraints? In this section we will propose a systematic method to address these  
547 questions in the context of the Holton-Mass model.

548 Consider a single-variable observable like  $U$  (30 km). If constrained to observe only  $U$  (30  
549 km) and forced to approximate  $q^+(x)$  as a function of this one variable, we would average  $q^+(\mathbf{x})$   
550 across the remaining 74 model dimensions, weighted by the invariant measure. We would assess  
551 the quality of this observable by the variance across those projected-out dimensions: a large  
552 projected variance would imply strong dependence on unobserved variables. Figure 5 applies  
553 this projection to the committor (first row) and mean first passage time to  $B$  (second row), using

554 three different single-variable observables:  $U$  (first column),  $\overline{v'T'}$  (second column), both at 30 km,  
 555 and the LASSO regression (third column). The solid curves show the projected means, and the  
 556 dotted curves indicate the one-standard-deviation envelope.  $q^+(U(30\text{ km}))$  is a smooth, mostly  
 557 monotonic curve with a consistently small projection error never exceeding  $\sim 0.2$ , which occurs  
 558 near  $q^+ = 0.5$ . Compare this to  $q^+(\overline{v'T'}(30\text{ km}))$ , which is essentially discontinuous as heat flux  
 559 increases from zero, and which has a large standard deviation approaching 0.3 when heat flux is  
 560 small. This is consistent with its prediction properties as shown in Figure 4: while the heat flux  
 561 reading hardly changes at all from zero, crucial processes are actively destabilizing the vortex, with  
 562 the committor increasing significantly without any response from  $\overline{v'T'}$ . The projected mean first  
 563 passage times tell a similar story, being strongly negatively correlated with the committor. Weaker  
 564 zonal wind generally signals less lead time before entering state  $B$ , but while heat flux stays small,  
 565 an observer is in the dark about how soon, as well as how certain, a transition is.

566 Let us briefly formalize the projection idea before exploring other variables. We want to approxi-  
 567 mate a function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ , such as the committor or mean first passage time, as a function of some  
 568 reduced coordinates  $\boldsymbol{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^k$ , called “collective variables” (CVs) in chemistry literature. That  
 569 is, we wish to find  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  such that  $F(\mathbf{x}) \approx f(\boldsymbol{\theta}(\mathbf{x}))$ . For instance,  $\boldsymbol{\theta}(\mathbf{x}) = (\theta_1(\mathbf{x}), \theta_2(\mathbf{x}))$  where  
 570  $\theta_1(\mathbf{x})$  is the mean zonal wind at 30 km and  $\theta_2(\mathbf{x})$  is the perturbation streamfunction magnitude  
 571  $|\Psi|$  at 30 km. Typically the projected dimension  $k \ll d$ , for instance  $k = 1$  or 2 for visualization  
 572 purposes. The “best” function  $f$  is chosen by minimizing some function-space metric between  
 573  $f \circ \boldsymbol{\theta}$  and  $F$ . The simplest choice would be the mean-squared error, so the projection problem is to  
 574 minimize over functions  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  the penalty

$$\begin{aligned}
 S[f; \boldsymbol{\theta}] &:= \|f \circ \boldsymbol{\theta} - F\|_{L^2(\pi)}^2 \\
 &= \int_{\mathbb{R}^d} \left[ f(\boldsymbol{\theta}(\mathbf{x})) - F(\mathbf{x}) \right]^2 \pi(d\mathbf{x})
 \end{aligned} \tag{20}$$

575 The optimal  $f$  for this purpose is the conditional expectation  $f(\mathbf{y}) = \mathbb{E}_{\mathbf{X} \sim \pi} [F(\mathbf{X}) | \boldsymbol{\theta}(\mathbf{X}) = \mathbf{y}] =$   
576  $\int f(\mathbf{x}) \delta(\boldsymbol{\theta}(\mathbf{x}) - \mathbf{y}) \pi(d\mathbf{x})$ . We derive a discretized version of this formula in the supplement, and  
577 this is how we display all the low-dimensional projections. We call the square root of  $S[f; \boldsymbol{\theta}]$  the  
578 projected standard deviation, or projection error, which determines the dotted envelope in Figure  
579 5.

580 A much harder problem than optimizing over  $f$  given  $\boldsymbol{\theta}$  is the problem of optimizing over sets of  
581 coordinates  $\boldsymbol{\theta}$ . CVs can be arbitrarily complex nonlinear functions of the basic state variables  $\mathbf{x}$ .  
582 Modern machine learning algorithms such as artificial neural networks are designed exactly for that  
583 purpose: to represent functions nonparametrically from observed input-output pairs. However, we  
584 wish to maintain some interpretability in the committor representation. For this reason, in searching  
585 for optimal projections, we begin with more constrained and physics-informed feature spaces before  
586 allowing for more complex relationships. We focus on observables coming from the Eliassen-Palm  
587 (EP) relation, which relates wave activity, PV fluxes and gradients, and heating source terms in a  
588 conservation equation. From Yoden (1987b), the EP relation for the Holton-Mass model takes the  
589 form

$$\begin{aligned} \partial_t \left( \frac{q'^2}{2} \right) + (\partial_y \bar{q}) \rho_s^{-1} \nabla \cdot \mathbf{F} \\ = - \frac{f_0^2}{N^2} \rho_s^{-1} \overline{q' \partial_z (\alpha \rho_s \partial_z \psi')} \end{aligned} \quad (21)$$

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

590 In the highly idealized Holton-Mass model, the EP flux divergence has two alternative expressions:  
591  $\rho_s^{-1} \nabla \cdot \mathbf{F} = \overline{v'q'} = \frac{R}{H f_0} \rho_s^{-1} \overline{v'T'}$ . If there were no dissipation ( $\alpha = 0$ ) and the background zonal state  
592 were time-independent ( $\partial_t \bar{q} = 0$ ), dividing both sides by  $\partial_y \bar{q}$  would express local conservation of  
593 wave activity  $\mathcal{A} = \overline{q'^2} / (2\partial_y \bar{q})$ . Neither of these is true in the stochastic Holton-Mass model, so  
594 we use the quantities in Equation (21) as diagnostics: enstrophy  $\overline{q'^2}$ , PV gradient  $\partial_y \bar{q}$ , PV flux  $\overline{v'q'}$ ,

595 and heat flux  $\overline{v'T'}$ . Each field is a function of  $(y, z)$  and takes on very different profiles in  $A$  and  $B$ ,  
596 as found by Yoden (1987b). A transition from  $A$  to  $B$ , where the vortex weakens dramatically, must  
597 entail a reduction in  $\partial_y \bar{q}$  and a burst in positive  $\overline{v'T'}$  and negative  $\overline{v'q'}$  as a Rossby wave propagates  
598 from the tropopause vertically up through the stratosphere. This is the general physical narrative  
599 of a sudden warming event, and these same fields might be expected to be useful observables to  
600 track for qualitative understanding and prediction, along with the basic state variables  $U$  and  $|\Psi|$ .

601 One option is to take vertical averages of any of these fields, but there may be particularly salient  
602 altitude levels that clarify the role of vertical interactions. The first three rows of Figure 6 display,  
603 for three of these fields ( $U$ ,  $|\Psi|$  and  $\overline{v'T'}$ ) and for a range of altitude levels, the mean and standard  
604 deviation of the committor projected onto that field at that altitude. Each altitude has a different  
605 range of the CV; for example, because  $U$  has a Dirichlet condition at the bottom and a Neumann  
606 condition at the top, the lower levels have a much smaller range of variability than the high levels.  
607 We also plot the integrated variance, or  $L^2$  projection error, at each level in the right-hand column.  
608 A low projected committor variance over  $U$  at altitude  $z_0$  means that the committor is mostly  
609 determined by the single observable  $U(z_0)$ , while a high projected variance indicates significant  
610 dependence of  $q^+$  on variables other than  $U(z_0)$ . In order to compare different altitudes and fields  
611 as directly as possible, the  $L^2$  projection error at each altitude is an average over discrete bins of  
612 the observable, not a proper integral.

613 In selecting good CV's, we generally look for a simple, hopefully monotonic, and sensitive  
614 relationship with the committor. Of all the candidate fields,  $U$  and  $\partial_y \bar{q}$  stand out the most in  
615 this respect, being clearly negatively correlated with the forward committor at all altitudes. The  
616 associated projection error tends to be greatest in the region  $q^+ \approx 0.5$ , as observed before, but  
617 interestingly there is a small altitude band around 20–25 km where its magnitude is minimized.  
618 This suggests an optimal altitude for monitoring the committor through zonal wind, giving the

619 most reliable estimate possible for a single state variable. In contrast, the projection of  $q^+$  onto  
620  $|\Psi|$ , displays a large variance across all altitudes. The eddy heat flux is also rather unhelpful as  
621 an early warning sign, despite its central role in SSW evolution, which is consistent with Figure  
622 5. For example, the large, positive spikes in heat flux across all altitudes generally occur after the  
623 committor  $\approx 0.5$  threshold has already been crossed. Furthermore, the relationship of  $\overline{v'T'}$  with  
624 the committor is not smooth. The  $q^+ < 0.5$  region at each altitude is a thin band near zero. Even  
625 so, the optimal altitude for observing the committor through heat flux is also 20 km.

626 The exhaustive observable search in Figure 6 is visually compelling, but not completely numer-  
627 ically satisfactory as a comparison between fields. Differences between units and ranges make it  
628 difficult to objectively compare the  $L^2$  projection error, despite the normalization mentioned above.  
629 Furthermore, restricting to one variable at a time is limiting. Accordingly, in a second, more auto-  
630 mated approach to identify salient variables, we perform a sparsity-promoting LASSO regression  
631 for the forward committor (Tibshirani 1996; Pedregosa et al. 2011), using as input features all state  
632 variables  $U, \text{Re}\Psi, \text{Im}\Psi$  and their vertical derivatives. We leave out eddy fluxes, which seem to  
633 have poor prediction properties. The advantage of a sparsity-promoting regression is to isolate a  
634 small number of observables that can decently approximate the committor in linear combination.  
635 Considering that regions close to  $A$  and  $B$  have low committor uncertainty, we regress only on data  
636 points with  $q^+ \in (0.2, 0.8)$ , and of those only a subset weighted by the reactive probability density  
637  $q^+ q^- \pi$ , since we wish to isolate the dominant transition pathways. To enforce predicted committors  
638 being between zero and one, we regress on the probit-transformed committor  $\ln(q^+/(1 - q^+))$ . First  
639 we do this at each altitude separately, and in Figure 7 (a) we plot the coefficients of each component  
640 as a function of altitude.

641 Each component is salient for some altitude range. In general,  $U$  and  $U_z$  dominate as causal  
642 variables at low altitudes, while  $\Psi$  and  $\Psi_z$  dominate at high altitudes. The overall prediction

643 quality, as measured by  $R^2$  and plotted in Figure 7 (b), is greatest around 21.5 km, consistent  
644 with our qualitative observations of Figure 6. Note that not all single-altitude slices are sufficient  
645 for approximating the committor, even with LASSO regression; in the altitude band 50 – 60 km,  
646 the LASSO predictor is not monotonic and has a large projected variance. The specific altitude  
647 can matter a great deal. But by using all altitudes at once, the committor approximation may be  
648 improved further. We thus repeat the LASSO with all altitudes simultaneously and find the sparse  
649 coefficient structure shown in 7 (c), with a few variables contributing the most:  $U$  (21.5 km),  
650  $U$  (29.6 km),  $\text{Re}\Psi_z$  (13.5 km), and  $\text{Im}\Psi$  (21.5 km). The results of LASSO regression are also  
651 displayed in the bottom row of Figure 4, the right column of Figure 5, and the bottom row of Figure  
652 6 for direct comparison with the other candidate fields. With multiple lines of evidence indicating  
653 21.5 km as an altitude with high predictive value for the forward committor, we can make a strong  
654 recommendation for targeting observations there. This conclusion applies only to the Holton-Mass  
655 model under these parameters, but the methodology explained above can be applied similarly to  
656 models of arbitrary complexity.

### 657 *c. Relationship to lead time*

658 A skillful forecast is only useful if it comes early and leaves some buffer time before impact.  
659 Having identified the committor as optimally skillful among all observables, we can now assess  
660 the limits of early prediction by relating certainty levels and lead times. Such a relationship would  
661 answer two dual questions: during the transition to an SSW winter phase, (1) how far in advance  
662 will we be aware of it with some prescribed confidence, say 80%? (2) given some prescribed lead  
663 time, say 42 days, how aware or in the dark could we be of it?

664 These questions clearly involve some kind of first-passage time, like the curves in the bottom  
665 row of Figure 5. The same quantity has been calculated previously in other simplified models, e.g.

666 Birner and Williams (2008) and Esler and Mester (2019). But  $\mathbb{E}[\tau_B]$  has an obvious shortcoming.  
 667 From Figure 5, we see that when  $q^+ \approx 0.5$ ,  $\mathbb{E}[\tau_B^+] \approx 600$ , an average which includes half the  
 668 paths going straight into  $B$  and the other half returning to  $A$  and lingering there before eventually  
 669 crossing into  $B$ . The conditional passage time  $\mathbb{E}[\tau_B | \tau_B < \tau_A]$  is designed to highlight only the  
 670 contribution of the latter half and measure the mean time of paths going directly to  $B$ , which can  
 671 be computed by DGA using a Laplace transform as described in Subsection 3(b). Figure 8 shows  
 672 all three quantities—the forward committor, mean passage time to  $B$ , and conditional passage time  
 673 to  $B$ —this time projected on a two-dimensional observable space  $(\text{Im}\Psi(21.5\text{km}), U(21.5\text{km}))$   
 674 identified as salient by sparse regression. Physically, these levels operate as a valve regulating wave  
 675 propagation into the stratosphere.

676 The committor has a clear negative relationship with both conditional and unconditional first  
 677 passage time: as the risk of imminent SSW grows, the time until impact shrinks. Figure 9 shows  
 678 this relationship more quantitatively, for both the  $A \rightarrow B$  process (panel (a)) and the  $B \rightarrow A$  process  
 679 (panel (b)). The relationship is roughly quantified by a least-squares regressions, weighted by the  
 680 change of measure, between the SSW probability  $q^+$  (resp. the restoration probability  $1 - q^+$ ) and  
 681 the conditional lead time to the SSW event  $\mathbb{E}[\tau_B | \tau_B < \tau_A]$  (resp. the conditional lead time to vortex  
 682 restoration,  $\mathbb{E}[\tau_A | \tau_A < \tau_B]$ ). While the relationships are nonlinear and the spreads significant, the  
 683 linear fits offer two meaningful numerical insight. The vertical intercept says how long the next  
 684 excursion to a given state will take when the system starts trapped in the other state. The negative  
 685 slope says how fast the remaining time shrinks as the risk grows. The vertical intercepts of 79 days  
 686 and 107 days offer further evidence that the vortex breaks down faster than it restores.

687 These metrics can inform preparation for extreme weather. For example, a threatened community  
 688 might decide in advance on an “alarm threshold” of, say, 50%, meaning they plan to prepare for an  
 689 SSW event only once it is 50% certain to occur. According to the linear fit in panel (a), they must

690 be ready to do so in  $\sim 48$  days time. The nonlinear deviations are, however, significant. The spread  
691 around the linear fit increases suddenly towards the lower-right corner of the plot, meaning that  
692 the uncertainty in timing, viewed as a function of the committor, increases as the SSW certainty  
693 increases. Lead time must therefore depend strongly on more than just the forward committor, and  
694 must be estimated by taking more details of the current state into account. We emphasize that the  
695 choice of  $A$ ,  $B$  and alarm thresholds are more of a community and policy decision than a scientific  
696 one. The strength of our approach is that it provides a flexible numerical framework to quantify  
697 and optimize the consequences of those decisions.

## 698 **6. Conclusion**

699 Forecasting rare events is, by the very nature of rare events, an extremely difficult computational  
700 task. Given the dangers posed by climate change, it also one of science's most pressing challenges.  
701 We suggest a computational framework that uses relatively short model simulations to make  
702 predictions on much longer time scales. Our numerical results point to its promise for forecasting.

703 Within the context of a stochastically forced Holton-Mass model with 75 degrees of freedom,  
704 we have computed fundamental quantities of the SSW transition process, including committor  
705 probabilities and expected lead times, for both the vortex destruction and vortex restoration pro-  
706 cesses. The system is irreversible, making these two directions very statistically distinct from each  
707 other. By systematically evaluating many model variables for their utility in predicting the fate  
708 of the vortex, we have identified some salient physical descriptions of early warning signs. We  
709 have furthermore quantified the relationship between probability and lead time for a given rare  
710 event, a potentially useful paradigm for assessing predictability and preparing for extreme weather.  
711 Our results suggest that the slow evolution of vortex preconditioning is an important source of

712 predictability. In particular, the zonal wind and streamfunction at 20 km seems to be optimal  
713 among a large class of dynamically motivated observables.

714 The committor and mean first passage time have obvious utility for forecasting, but they are also  
715 ingredients in a larger framework called Transition Path Theory (TPT) for describing rare steady  
716 state transition events. In principle, interrogating the ensemble of transition paths requires direct  
717 simulation of the system long enough to observe many transition events. However, using TPT,  
718 quantities computable by our framework can be combined to yield key statistics describing the  
719 ensemble of *transition paths* connecting regions in state space, (Finkel et al. 2020; Metzner et al.  
720 2006, 2009; Vanden-Eijnden and E 2010; E. and Vanden-Eijnden 2006). In a following paper we  
721 will apply the same short-trajectory forecasting approach together with TPT to compute transition  
722 path statistics such as return times and extract insight about physical mechanisms of the transition  
723 process.

724 Our numerical pipeline is promising and robust, but leaves plenty of room for improvement.  
725 Our sampling method, while advantageous for validation of results, wastes a great deal of data.  
726 Targeted sampling from the transition region has the potential to achieve the same precision for the  
727 quantities of interest with much less data. Also, moving beyond a basis expansion of the forecast  
728 functions, in upcoming work we will explore more flexible representations using kernel methods  
729 and neural networks. The solution of high-dimensional PDEs is an active research area that is  
730 making innovative use of machine learning, particularly in the fields of computational chemistry  
731 and quantum mechanics (e.g., Chen and Majda 2017; Carleo and Troyer 2017; Han et al. 2018;  
732 Khoo et al. 2018; Li et al. 2020; Mardt et al. 2018; Li et al. 2019; Lorpaiboon et al. 2020). Similar  
733 approaches may hold great potential for understanding predictability in atmospheric science.

734 *Acknowledgments.* J.F. is supported by the U.S. Department of Energy, Office of Science, Office of  
735 Advanced Scientific Computing Research, Department of Energy Computational Science Graduate  
736 Fellowship under Award Number DE-SC0019323. R.J.W. is supported by New York University's  
737 Dean's Dissertation Fellowship and by the Research Training Group in Modeling and Simulation  
738 funded by the National Science Foundation via grant RTG/DMS-1646339. E.P.G. acknowledges  
739 support from the U. S. National Science Foundation through grant AGS-1852727. This work was  
740 partially supported by the NASA Astrobiology Program, grant No. 80NSSC18K0829 and benefited  
741 from participation in the NASA Nexus for Exoplanet Systems Science research coordination  
742 network. J.W. acknowledges support from the Advanced Scientific Computing Research Program  
743 within the DOE Office of Science through award DE-SC0020427. The computations in the  
744 paper were done on the high-performance computing clusters at New York University and the  
745 Research Computing Center at the University of Chicago. We thank John Strahan, Aaron Dinner,  
746 and Chatipat Lorpaiboon for helpful methodological advice. Mary Silber, Noboru Nakamura,  
747 and Richard Kleeman offered invaluable scientific insight. J.F. benefitted from many helpful  
748 discussions with Anya Katsevich.

749 *Data availability statement.* The stochastic Holton-Mass model and analysis techniques are fully  
750 described in the text and references, and can be integrated quickly at low computational costs. J.F.  
751 is happy to share code and data upon request.

## 752 **References**

753 Berry, T., D. Giannakis, and J. Harlim, 2015: Nonparametric forecasting of low-dimensional  
754 dynamical systems. *Phys. Rev. E*, **91**, 032915, doi:10.1103/PhysRevE.91.032915.

- 755 Birner, T., and P. D. Williams, 2008: Sudden stratospheric warmings as noise-induced transitions.  
756 *Journal of the Atmospheric Sciences*, **65** (10), 3337–3343, doi:10.1175/2008JAS2770.1.
- 757 Bolton, T., and L. Zanna, 2019: Applications of deep learning to ocean data in-  
758 ference and subgrid parameterization. *Journal of Advances in Modeling Earth Sys-*  
759 *tems*, **11** (1), 376–399, doi:<https://doi.org/10.1029/2018MS001472>, URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001472>, <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018MS001472>.
- 762 Bouchet, F., J. Laurie, and O. Zaboronski, 2011: Control and instanton trajectories for ran-  
763 dom transitions in turbulent flows. *Journal of Physics: Conference Series*, **318** (2), 022 041,  
764 doi:10.1088/1742-6596/318/2/022041, URL <https://doi.org/10.1088/1742-6596/318/2/022041>.
- 766 Bouchet, F., J. Laurie, and O. Zaboronski, 2014: Langevin dynamics, large deviations and instan-  
767 tons for the quasi-geostrophic model and two-dimensional euler equations. *Journal of Statis-*  
768 *tical Physics*, **156**, 1066–1092, doi:10.1007/s10955-014-1052-5, URL <https://doi.org/10.1007/s10955-014-1052-5>.
- 770 Bouchet, F., J. Rolland, and E. Simonnet, 2019a: Rare event algorithm links transitions in  
771 turbulent flows with activated nucleations. *Physical Review Letters*, **122** (7), 074 502, doi:  
772 [10.1103/PhysRevLett.122.074502](https://doi.org/10.1103/PhysRevLett.122.074502).
- 773 Bouchet, F., J. Rolland, and J. Wouters, 2019b: Rare event sampling methods. *Chaos: An Inter-*  
774 *disciplinary Journal of Nonlinear Science*, **29** (8), 080 402, doi:10.1063/1.5120509.
- 775 Bowman, G. R., V. S. Pande, and F. Noé, 2013: *An introduction to Markov state models and*  
776 *their application to long timescale molecular simulation*, Vol. 797. Springer Science & Business

777 Media.

778 Butler, A. H., D. J. Seidel, S. C. Hardiman, N. Butchart, T. Birner, and A. Match, 2015: Defining  
779 sudden stratospheric warmings. *Bulletin of the American Meteorological Society*, **96** (11), 1913–  
780 1928, doi:10.1175/BAMS-D-13-00173.1.

781 Carleo, G., and M. Troyer, 2017: Solving the quantum many-body problem with arti-  
782 ficial neural networks. *Science*, **355** (6325), 602–606, doi:10.1126/science.aag2302, URL  
783 <https://science.sciencemag.org/content/355/6325/602>, [https://science.sciencemag.org/content/](https://science.sciencemag.org/content/355/6325/602.full.pdf)  
784 [355/6325/602.full.pdf](https://science.sciencemag.org/content/355/6325/602.full.pdf).

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

788 Charney, J. G., and J. G. DeVore, 1979: Multiple Flow Equilibria in the At-  
789 mosphere and Blocking. *Journal of the Atmospheric Sciences*, **36** (7), 1205–1216,  
790 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)  
791 [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_1205_mfeita_2_0_co_2.pdf)  
792 [36/7/1205/3420739/1520-0469\(1979\)036\\_1205\\_mfeita\\_2\\_0\\_co\\_2.pdf](https://journals.ametsoc.org/jas/article-pdf/36/7/1205/3420739/1520-0469(1979)036_1205_mfeita_2_0_co_2.pdf).

793 Charney, J. G., and P. G. Drazin, 1961: Propagation of planetary-scale disturbances from  
794 the lower into the upper atmosphere. *Journal of Geophysical Research (1896-1977)*,  
795 **66** (1), 83–109, doi:10.1029/JZ066i001p00083, URL [https://agupubs.onlinelibrary.wiley.com/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JZ066i001p00083)  
796 [doi/abs/10.1029/JZ066i001p00083](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JZ066i001p00083), [https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/JZ066i001p00083)  
797 [JZ066i001p00083](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/JZ066i001p00083).

798 Chattopadhyay, A., E. Nabizadeh, and P. Hassanzadeh, 2020: Analog forecast-  
799 ing of extreme-causing weather patterns using deep learning. *Journal of Ad-  
800 vances in Modeling Earth Systems*, **12** (2), e2019MS001958, doi:[https://doi.  
801 org/10.1029/2019MS001958](https://doi.org/10.1029/2019MS001958), URL [https://agupubs.onlinelibrary.  
802 2019MS001958](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001958), e2019MS001958 10.1029/2019MS001958, [https://agupubs.onlinelibrary.  
803 wiley.com/doi/pdf/10.1029/2019MS001958](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2019MS001958).

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

808 Chen, N., and A. J. Majda, 2017: Beating the curse of dimension with accurate statistics for the  
809 fokker–planck equation in complex turbulent systems. *Proceedings of the National Academy of  
810 Sciences*, **114** (49), 12 864–12 869, doi:10.1073/pnas.1717017114, URL [https://www.pnas.org/  
811 content/114/49/12864](https://www.pnas.org/content/114/49/12864), <https://www.pnas.org/content/114/49/12864.full.pdf>.

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

816 Chodera, J. D., and F. Noé, 2014: Markov state models of biomolecular conformational dynamics.  
817 *Current Opinion in Structural Biology*, **25**, 135 – 144, doi:[https://doi.org/10.1016/j.sbi.2014.04.  
818 002](https://doi.org/10.1016/j.sbi.2014.04.002), URL <http://www.sciencedirect.com/science/article/pii/S0959440X14000426>, theory and  
819 simulation / Macromolecular machines.

- 820 Christiansen, B., 2000: Chaos, quasiperiodicity, and interannual variability: Studies of a strato-  
821 spheric vacillation model. *Journal of the Atmospheric Sciences*, **57** (18), 3161–3173, doi:  
822 10.1175/1520-0469(2000)057<3161:CQAIVS>2.0.CO;2.
- 823 Crommelin, D. T., 2003: Regime transitions and heteroclinic connections in a barotropic  
824 atmosphere. *Journal of the Atmospheric Sciences*, **60** (2), 229 – 246, doi:10.  
825 1175/1520-0469(2003)060<0229:RTAHCI>2.0.CO;2, URL [https://journals.ametsoc.org/view/  
826 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).
- 827 DelSole, T., and B. F. Farrell, 1995: A stochastically excited linear system as a model for quasi-  
828 geostrophic turbulence: Analytic results for one- and two-layer fluids. *Journal of the Atmospheric  
829 Sciences*, **52** (14), 2531–2547, doi:10.1175/1520-0469(1995)052<2531:ASELSA>2.0.CO;2.
- 830 Dematteis, G., T. Grafke, M. Onorato, and E. Vanden-Eijnden, 2019: Experimental evidence  
831 of hydrodynamic instantons: The universal route to rogue waves. *Phys. Rev. X*, **9**, 041057,  
832 doi:10.1103/PhysRevX.9.041057, URL <https://link.aps.org/doi/10.1103/PhysRevX.9.041057>.
- 833 Dematteis, G., T. Grafke, and E. Vanden-Eijnden, 2018: Rogue waves and large deviations in deep  
834 sea. *Proceedings of the National Academy of Sciences*, **115** (5), 855–860, doi:10.1073/pnas.  
835 1710670115.
- 836 Durrett, R., 2013: *Probability: Theory and Examples*. Cambridge University Press.
- 837 E., W., and E. Vanden-Eijnden, 2006: Towards a Theory of Transition Paths. *Journal of Sta-  
838 tistical Physics*, **123** (3), 503, doi:10.1007/s10955-005-9003-9, URL [https://doi.org/10.1007/  
839 s10955-005-9003-9](https://doi.org/10.1007/s10955-005-9003-9).
- 840 Esler, J. G., and M. Mester, 2019: Noise-induced vortex-splitting stratospheric sudden warmings.  
841 *Quarterly Journal of the Royal Meteorological Society*, **145** (719), 476–494, doi:<https://doi.org/>

842 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)  
843 [rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.3443](https://rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.3443).

844 Farazmand, M., and T. P. Sapsis, 2017: A variational approach to probing extreme events in turbu-  
845 lent dynamical systems. *Science Advances*, **3** (9), doi:10.1126/sciadv.1701533, URL [https://](https://advances.sciencemag.org/content/3/9/e1701533)  
846 [advances.sciencemag.org/content/3/9/e1701533](https://advances.sciencemag.org/content/3/9/e1701533), [https://advances.sciencemag.org/content/3/9/](https://advances.sciencemag.org/content/3/9/e1701533.full.pdf)  
847 [e1701533.full.pdf](https://advances.sciencemag.org/content/3/9/e1701533.full.pdf).

848 Finkel, J., D. S. Abbot, and J. Weare, 2020: Path Properties of Atmospheric Transitions: Illustra-  
849 tion with a Low-Order Sudden Stratospheric Warming Model. *Journal of the Atmospheric*  
850 *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)  
851 [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)  
852 [pdf](https://journals.ametsoc.org/jas/article-pdf/77/7/2327/4958190/jasd190278.pdf).

853 Frank, N., and S. Fischer, 2008: Transition networks for modeling the kinetics of conformational  
854 change in macromolecules. *Current Opinion in Structural Biology*, **18**, 154–163, doi:10.1016/  
855 [j.sbi.2008.01.008](https://doi.org/10.1016/j.sbi.2008.01.008).

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

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

861 Gottwald, G. A., D. T. Crommelin, and C. L. E. Franzke, 2016: Stochastic climate theory.  
862 1612.07474.

- 863 Han, J., A. Jentzen, and W. E, 2018: Solving high-dimensional partial differential equations  
864 using deep learning. *Proceedings of the National Academy of Sciences*, **115 (34)**, 8505–8510,  
865 doi:10.1073/pnas.1718942115, URL <https://www.pnas.org/content/115/34/8505>, <https://www.pnas.org/content/115/34/8505.full.pdf>.  
866
- 867 Hasselmann, K., 1976: Stochastic climate models part i. theory. *Tellus*, **28 (6)**, 473–485, doi:  
868 10.3402/tellusa.v28i6.11316.
- 869 Hoffman, R. N., J. M. Henderson, S. M. Leidner, C. Grassotti, and T. Nehr Korn, 2006: The response  
870 of damaging winds of a simulated tropical cyclone to finite-amplitude perturbations of different  
871 variables. *Journal of the Atmospheric Sciences*, **63 (7)**, 1924 – 1937, doi:10.1175/JAS3720.1,  
872 URL <https://journals.ametsoc.org/view/journals/atsc/63/7/jas3720.1.xml>.
- 873 Holton, J. R., and C. Mass, 1976: Stratospheric vacillation cycles. *Journal of the Atmospheric*  
874 *Sciences*, **33 (11)**, 2218–2225, doi:10.1175/1520-0469(1976)033<2218:SVC>2.0.CO;2.
- 875 Hu, G., T. Bódai, and V. Lucarini, 2019: Effects of stochastic parametrization on extreme value  
876 statistics. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, **29 (8)**, 083 102, doi:10.  
877 1063/1.5095756, URL <https://doi.org/10.1063/1.5095756>, <https://doi.org/10.1063/1.5095756>.
- 878 Karatzas, I., and S. E. Shreve, 1998: *Brownian Motion and Stochastic Calculus*. Springer.
- 879 Khoo, Y., J. Lu, and L. Ying, 2018: Solving for high-dimensional committor functions  
880 using artificial neural networks. *Research in the Mathematical Sciences*, **6**, doi:10.1007/  
881 s40687-018-0160-2, URL <https://doi.org/10.1007/s40687-018-0160-2>.
- 882 L'Ecuyer, P., V. Demers, and B. Tuffin, 2007: Rare events, splitting, and quasi-monte carlo. *ACM*  
883 *Transactions on Modeling and Computer Simulation (TOMACS)*, **17 (2)**, 9–es.

- 884 Legras, B., and M. Ghil, 1985: Persistent anomalies, blocking and variations in at-  
885 mospheric predictability. *Journal of Atmospheric Sciences*, **42** (5), 433 – 471, doi:10.  
886 1175/1520-0469(1985)042<0433:PABAVI>2.0.CO;2, URL [https://journals.ametsoc.org/view/  
887 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).
- 888 Li, H., Y. Khoo, Y. Ren, and L. Ying, 2020: Solving for high dimensional committor functions  
889 using neural network with online approximation to derivatives. 2012.06727.
- 890 Li, Q., B. Lin, and W. Ren, 2019: Computing committor functions for the study of rare events using  
891 deep learning. *The Journal of Chemical Physics*, **151** (5), 054 112, doi:10.1063/1.5110439, URL  
892 <https://doi.org/10.1063/1.5110439>.
- 893 Lorenz, E. N., 1963: Deterministic nonperiodic flow. *Journal of Atmospheric Sciences*, **20** (2), 130  
894 – 141, doi:10.1175/1520-0469(1963)020<0130:DNF>2.0.CO;2, URL [https://journals.ametsoc.  
895 org/view/journals/atsc/20/2/1520-0469\\_1963](https://journals.ametsoc.org/view/journals/atsc/20/2/1520-0469_1963).
- 896 Lorpaiboon, C., E. H. Thiede, R. J. Webber, J. Weare, and A. R. Dinner, 2020: Integrated  
897 variational approach to conformational dynamics: A robust strategy for identifying eigenfunc-  
898 tions of dynamical operators. *The Journal of Physical Chemistry B*, **124** (42), 9354–9364,  
899 doi:10.1021/acs.jpcc.0c06477, URL <https://doi.org/10.1021/acs.jpcc.0c06477>.
- 900 Lucente, D., S. Duffner, C. Herbert, J. Rolland, and F. Bouchet, 2019: Machine learning of  
901 committor functions for predicting high impact climate events. *Climate Informatics*, Paris,  
902 France, URL <https://hal.archives-ouvertes.fr/hal-02322370>.
- 903 Majda, A. J., and D. Qi, 2018: Strategies for reduced-order models for predicting the statistical  
904 responses and uncertainty quantification in complex turbulent dynamical systems. *SIAM Review*,

905 **60 (3)**, 491–549, doi:10.1137/16M1104664, URL <https://doi.org/10.1137/16M1104664>, <https://doi.org/10.1137/16M1104664>.  
906

907 Majda, A. J., I. Timofeyev, and E. Vanden Eijnden, 2001: A mathematical framework for  
908 stochastic climate models. *Communications on Pure and Applied Mathematics*, **54 (8)**,  
909 891–974, doi:<https://doi.org/10.1002/cpa.1014>, URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpa.1014>, <https://onlinelibrary.wiley.com/doi/pdf/10.1002/cpa.1014>.  
910

911 Mardt, A., L. Pasuali, H. Wu, and F. Noé, 2018: Vampnets for deep learning of molecular kinetics.  
912 *Nature Communications*, **9**, doi:10.1038/s41467-017-02388-1, URL <https://doi.org/10.1038/s41467-017-02388-1>.  
913

914 Matsuno, T., 1971: A Dynamical Model of the Stratospheric Sudden Warming. *Journal*  
915 *of the Atmospheric Sciences*, **28 (8)**, 1479–1494, doi:10.1175/1520-0469(1971)028<1479:  
916 ADMOTS>2.0.CO;2, URL [https://doi.org/10.1175/1520-0469\(1971\)028<1479:ADMOTS>2.0.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%5B1479%5D\\_admots%5B2%5D\\_co%5B2%5D.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).  
918

919 Metzner, P., C. Schutte, and E. Vanden-Eijnden, 2006: Illustration of transition path theory  
920 on a collection of simple examples. *The Journal of Chemical Physics*, **125 (8)**, 1–2, doi:  
921 10.1063/1.2335447.

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

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

- 927 Mohamad, M. A., and T. P. Sapsis, 2018: Sequential sampling strategy for extreme event  
928 statistics in nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*,  
929 **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)  
930 [115/44/11138](https://www.pnas.org/content/115/44/11138), <https://www.pnas.org/content/115/44/11138.full.pdf>.
- 931 Ngwira, C. M., and Coauthors, 2013: Simulation of the 23 July 2012 extreme  
932 space weather event: What if this extremely rare cme was earth directed? *Space*  
933 *Weather*, **11 (12)**, 671–679, doi:<https://doi.org/10.1002/2013SW000990>, URL [https://agupubs.](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2013SW000990)  
934 [onlinelibrary.wiley.com/doi/abs/10.1002/2013SW000990](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2013SW000990), [https://agupubs.onlinelibrary.wiley.](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1002/2013SW000990)  
935 [com/doi/pdf/10.1002/2013SW000990](https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1002/2013SW000990).
- 936 Noé, F., C. Schütte, E. Vanden-Eijnden, L. Reich, and T. R. Weigl, 2009: Constructing the equilib-  
937 rium ensemble of folding pathways from short off-equilibrium simulations. *Proceedings of the*  
938 *National Academy of Sciences*, **106 (45)**, 19 011–19 016, doi:10.1073/pnas.0905466106, URL  
939 <https://www.pnas.org/content/106/45/19011>, [https://www.pnas.org/content/106/45/19011.full.](https://www.pnas.org/content/106/45/19011.full.pdf)  
940 [pdf](https://www.pnas.org/content/106/45/19011.full.pdf).
- 941 Oksendal, B., 2003: *Stochastic Differential Equations: An Introduction with Applications*.  
942 Springer.
- 943 Pande, V. S., K. Beauchamp, and G. R. Bowman, 2010: Everything you wanted to know about  
944 markov state models but were afraid to ask. *Methods*, **52 (1)**, 99–105, URL [https://doi.org/10.](https://doi.org/10.1016/j.ymeth.2010.06.002)  
945 [1016/j.ymeth.2010.06.002](https://doi.org/10.1016/j.ymeth.2010.06.002).
- 946 Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine learning in Python. *Journal of Machine*  
947 *Learning Research*, **12**, 2825–2830.

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

951 Ragone, F., and F. Bouchet, 2020: Computation of extreme values of time averaged observables  
952 in climate models with large deviation techniques. *Journal of Statistical Physics*, **179** (5), 1637–  
953 1665, doi:10.1007/s10955-019-02429-7, URL <https://doi.org/10.1007/s10955-019-02429-7>.

954 Ragone, F., J. Wouters, and F. Bouchet, 2018: Computation of extreme heat waves in climate  
955 models using a large deviation algorithm. *Proceedings of the National Academy of Sciences*,  
956 **115** (1), 24–29, doi:10.1073/pnas.1712645115, <https://www.pnas.org/content/115/1/24.full.pdf>.

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

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

962 Sapsis, T. P., 2021: Statistics of extreme events in fluid flows and waves. *Annual Review of Fluid*  
963 *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)  
964 [1146/annurev-fluid-030420-032810](https://doi.org/10.1146/annurev-fluid-030420-032810), <https://doi.org/10.1146/annurev-fluid-030420-032810>.

965 Schaller, N., J. Sillmann, J. Anstey, E. M. Fischer, C. M. Grams, and S. Russo, 2018: Influence of  
966 blocking on northern european and western russian heatwaves in large climate model ensembles.  
967 *Environmental Research Letters*, **13** (5), 054 015, doi:10.1088/1748-9326/aaba55, URL [https:](https://doi.org/10.1088/1748-9326/aaba55)  
968 [//doi.org/10.1088/1748-9326/aaba55](https://doi.org/10.1088/1748-9326/aaba55).

969 Sillmann, J., and Coauthors, 2017: Understanding, modeling and predicting weather and climate  
970 extremes: Challenges and opportunities. *Weather and Climate Extremes*, **18**, 65 – 74, doi:  
971 <https://doi.org/10.1016/j.wace.2017.10.003>.

972 Simonnet, E., J. Rolland, and F. Bouchet, 2020: Multistability and rare spontaneous transitions be-  
973 tween climate and jet configurations in a barotropic model of the jovian mid-latitude troposphere.  
974 2009.09913.

975 Strahan, J., A. Antoszewski, C. Lorpaiboon, B. P. Vani, J. Weare, and A. R. Dinner, 2020:  
976 Long-timescale predictions from short-trajectory data: A benchmark analysis of the trp-cage  
977 miniprotein. 2009.04034.

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

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

984 Tibshirani, R., 1996: Regression shrinkage and selection via the lasso. *Journal of the Royal*  
985 *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)  
986 [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)  
987 [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)  
988 [1996.tb02080.x](https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1996.tb02080.x).

989 Timmermann, A., F.-F. Jin, and J. Abshagen, 2003: A nonlinear theory for el niño bursting. *Jour-*  
990 *nal of the Atmospheric Sciences*, **60** (1), 152 – 165, doi:10.1175/1520-0469(2003)060<0152:

991 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_)  
992 2003.

993 Vanden-Eijnden, E., and W. E, 2010: Transition-path theory and path-finding algorithms for the  
994 study of rare events. *Annual Review of Physical Chemistry*, **61** (1), 391–420, doi:10.1146/  
995 annurev.physchem.040808.090412.

996 Vanden-Eijnden, E., and J. Weare, 2013: Data assimilation in the low noise regime with application  
997 to the kuroshio. *Monthly Weather Review*, **141** (6), 1822–1841, doi:10.1175/MWR-D-12-00060.  
998 1.

999 Vitart, F., and A. W. Robertson, 2018: The sub-seasonal to seasonal prediction project (s2s)  
1000 and the prediction of extreme events. *npj Climate and Atmospheric Science*, **1**, URL <https://doi.org/10.1038/s41612-018-0013-0>.  
1001

1002 Wan, Z. Y., P. Vlachas, P. Koumoutsakos, and T. Sapsis, 2018: Data-assisted reduced-order  
1003 modeling of extreme events in complex dynamical systems. *PLOS ONE*, **13** (5), 1–22, doi:  
1004 10.1371/journal.pone.0197704, URL <https://doi.org/10.1371/journal.pone.0197704>.

1005 Weare, J., 2009: Particle filtering with path sampling and an application to a bimodal ocean current  
1006 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)  
1007 [jcp.2009.02.033](https://doi.org/10.1016/j.jcp.2009.02.033).

1008 Webber, R. J., D. A. Plotkin, M. E. O’Neill, D. S. Abbot, and J. Weare, 2019: Practical rare event  
1009 sampling for extreme mesoscale weather. *Chaos*, **29** (5), 053 109, doi:10.1063/1.5081461.

1010 Weinan, E., T. Li, and E. Vanden-Eijnden, 2019: *Applied stochastic analysis*, Vol. 199. American  
1011 Mathematical Soc.

- 1012 Yasuda, Y., F. Bouchet, and A. Venaille, 2017: A new interpretation of vortex-split sudden  
1013 stratospheric warmings in terms of equilibrium statistical mechanics. *Journal of the Atmospheric*  
1014 *Sciences*, **74** (12), 3915–3936, doi:10.1175/JAS-D-17-0045.1.
- 1015 Yoden, S., 1987a: Bifurcation properties of a stratospheric vacillation model. *Journal of the At-*  
1016 *mospheric Sciences*, **44** (13), 1723–1733, doi:10.1175/1520-0469(1987)044<1723:BPOASV>  
1017 2.0.CO;2.
- 1018 Yoden, S., 1987b: Dynamical Aspects of Stratospheric Vacillations in a Highly  
1019 Truncated Model. *Journal of the Atmospheric Sciences*, **44** (24), 3683–3695,  
1020 doi:10.1175/1520-0469(1987)044<3683:DAOSVI>2.0.CO;2, URL [https://doi.org/10.1175/  
1021 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).
- 1022 Zhang, F., and J. A. Sippel, 2009: Effects of moist convection on hurricane predictability. *Journal*  
1023 *of the Atmospheric Sciences*, **66** (7), 1944 – 1961, doi:10.1175/2009JAS2824.1, URL [https:  
1024 //journals.ametsoc.org/view/journals/atasc/66/7/2009jas2824.1.xml](https://journals.ametsoc.org/view/journals/atasc/66/7/2009jas2824.1.xml).
- 1025 Zwanzig, R., 2001: *Nonequilibrium statistical mechanics*. Oxford University Press.

1026 **LIST OF TABLES**

1027 **Table 1.** Fraction of time that the system spends (i) inside  $A$ , (ii) outside  $A$  but destined  
1028 for  $A$ , (iii) inside  $B$ , (iv) outside  $B$  but destined for  $B$ . Each fraction is computed  
1029 directly from DGA and empirically verified from the long simulation. . . . . 51

Region	Fraction (DGA)	Fraction (empirical)
Inside $A$	0.24	0.27
$(A \cup B)^c \rightarrow A$	0.27	0.28
Inside $B$	0.14	0.13
$(A \cup B)^c \rightarrow B$	0.36	0.32

1030 TABLE 1. Fraction of time that the system spends (i) inside  $A$ , (ii) outside  $A$  but destined for  $A$ , (iii) inside  $B$ ,  
1031 (iv) outside  $B$  but destined for  $B$ . Each fraction is computed directly from DGA and empirically verified from  
1032 the long simulation.

## LIST OF FIGURES

1033		
1034	<b>Fig. 1.</b>	<b>Illustration of the two stable states of the Holton-Mass model and transitions between them.</b> (a) Zonal wind profiles of the radiatively maintained strong vortex (the fixed point <b>a</b> , blue) which increases linearly with altitude, and the weak vortex (the fixed point <b>b</b> , red) which dips close to zero in the mid-stratosphere. (b) Streamfunction contours are overlaid for the two equilibria <b>a</b> and <b>b</b> , the weak vortex exhibiting strong westward phase tilt with altitude. (c) Timeseries of $U(30\text{ km})$ from a long stochastic simulation, including several noise-induced transitions from <i>A</i> to <i>B</i> (green) and from <i>B</i> to <i>A</i> (orange). Although both states <b>a</b> and <b>b</b> are equilibria in this parameter regime ( $h = 38.5\text{ m}$ ), the stochastic perturbations uncover the vacillation cycles that would appear beyond the Hopf bifurcation if $h$ were increased. (d) A parametric curve of the same trajectory segment as in (c) with the same color highlights for transition paths, but in the space $( \Psi , U)$ at 30 km. The two equilibria are indicated with horizontal blue and red lines. . . . . 54
1035		
1036		
1037		
1038		
1039		
1040		
1041		
1042		
1043		
1044		
1045		
1046	<b>Fig. 2.</b>	<b>Accuracy of the committor and mean first passage time calculations verified with long trajectory data.</b> The DGA calculations assign approximate committor and mean first passage time values $q^+(\mathbf{X}_n(0))$ and $m_B(\mathbf{X}_n(0))$ to each data point. Because each snapshot $x_n$ was collected from a long trajectory, its destination and the time to get there are known and can provide an empirical validation of the committor and lead time. For 20 equal partitions $(\zeta_1, \zeta_2)$ of the interval $(0, 1)$ , we assemble all trajectory starts $\mathbf{X}_n(0)$ with $q^+(\mathbf{X}_n(0)) \in (\zeta_1, \zeta_2)$ and count the fraction heading toward <i>B</i> . These are the empirical committors for the interval $(\zeta_1, \zeta_2)$ , and are plotted on the vertical axis against $(\zeta_1 + \zeta_2)/2$ . Similarly, we bin the space of calculated first passage times and for each bin average the empirical first passage time to <i>B</i> . Both quantities line up well between DGA computations and empirical values, with the exception of the longest passage times, which are underestimated. . . . . 55
1047		
1048		
1049		
1050		
1051		
1052		
1053		
1054		
1055		
1056		
1057	<b>Fig. 3.</b>	<b>Steady-state distribution.</b> The density $\pi(\mathbf{x})$ is projected onto the two-dimensional space $( \Psi , U)$ at 30 km, on a log scale. The density is peaked in the neighborhoods of the two fixed points. On the right is a projection of $\pi$ onto the single variable $U(30\text{ km})$ , on a linear scale, confirming strong bimodality. . . . . 56
1058		
1059		
1060		
1061	<b>Fig. 4.</b>	<b>The committor vs. other observables as a forecasting tool.</b> A representative simulated SSW event from the long simulation is plotted over time, starting 65 days in advance of the official event when $U(30\text{ km})$ first drops below 1.75 m/s, which is marked by a vertical solid line. Panel (a) shows the committor over time following the trajectory, panel (b) shows the zonal wind $U(30\text{ km})$ , and panel (c) shows the eddy heat flux $\overline{v'T'}(30\text{ km})$ . Horizontal dashed lines mark the natural forecasting threshold of $q^+ = 0.5$ (panel (a)) or the value of the observable most closely associated with $q^+ = 0.5$ : 37 m/s (panel (b)) and $1.2 \times 10^{-6} \text{ K} \cdot \text{m/s}$ (panel (c)). The sharp increase in $q^+$ as it crosses the threshold provides a clear and early warning sign of oncoming SSW, about 26 days in advance. $U$ and $\overline{v'T'}$ are moving slowly at that time, and don't clear their respective thresholds for the last time until the event is much closer at hand. The gap in lead time is marked by blue strips. . . . . 57
1062		
1063		
1064		
1065		
1066		
1067		
1068		
1069		
1070		
1071		
1072	<b>Fig. 5.</b>	<b>One-dimensional projections of the forward committor and mean first passage time to <i>B</i>, computed with DGA.</b> These functions depend on all $d = 75$ degrees of freedom in the model, but we have averaged across $d - 1 = 74$ dimensions to visualize the committor (first row) and mean first passage time to <i>B</i> (second row) as rough functions of three single degrees of freedom: $U(30\text{ km})$ (first column), $\overline{v'T'}(30\text{ km})$ (second column), and the LASSO-regressed committor (third column). The forward committor measures proximity to <i>B</i> in probability, while mean passage time to <i>B</i> measures proximity in time, hence the negative correlation between the two quantities. The general trends reveal fairly obvious relationships: stronger wind is associated with tendency towards the strong-vortex state <i>A</i> , and larger poleward eddy
1073		
1074		
1075		
1076		
1077		
1078		
1079		
1080		

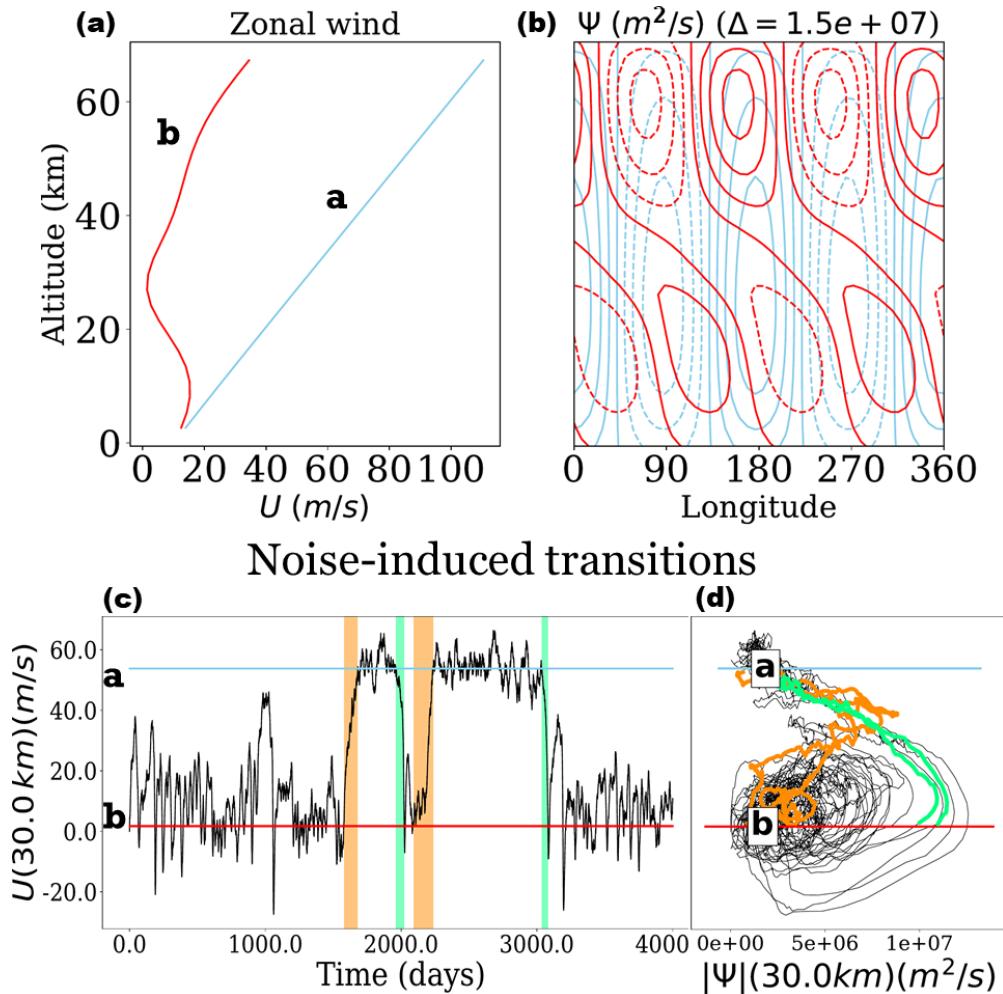
1081 heat flux is associated with tendency toward the weak vortex state  $B$ . In addition, curves like  
 1082 this assess the quality of single-variable observables as proxies for an oncoming transition  
 1083 event. The committor and passage time vary smoothly and (mostly) monotonically with  $U$ ,  
 1084 but discontinuously with  $\overline{v'T'}$ : the heat flux burst that accompanies a SSW gives no advance  
 1085 warning for the event, while a small negative change in  $U$  indicates incrementally higher  
 1086 transition probability and shorter lead time. . . . . 58

1087 **Fig. 6. Projection of the forward committor onto a large collection of one-dimensional CVs,**  
 1088 along with the associated standard deviation, or projection error, of the committor along the  
 1089 remaining 74 model dimensions. Consider the first two panels. The left-hand panel shows,  
 1090 for each discretized altitude  $z$ , a heatmap of the committor as  $U(z)$  ranges from its minimum  
 1091 to its maximum realized strength at that altitude. At the bottom is an additional heatmap of  
 1092 the committor as a function of  $\langle U \rangle_z$ , the vertical average. These are conditional expectations,  
 1093 with the corresponding conditional standard deviations displayed in the right-hand panels.  
 1094 The following rows display analogous plots for wave magnitude, eddy PV flux, background  
 1095 PV gradient, eddy heat flux, and the LASSO predictor specified in Figure 7. The bottom  
 1096 line on the last plot is not a vertical average, but the results of regression on all altitudes at once. 59

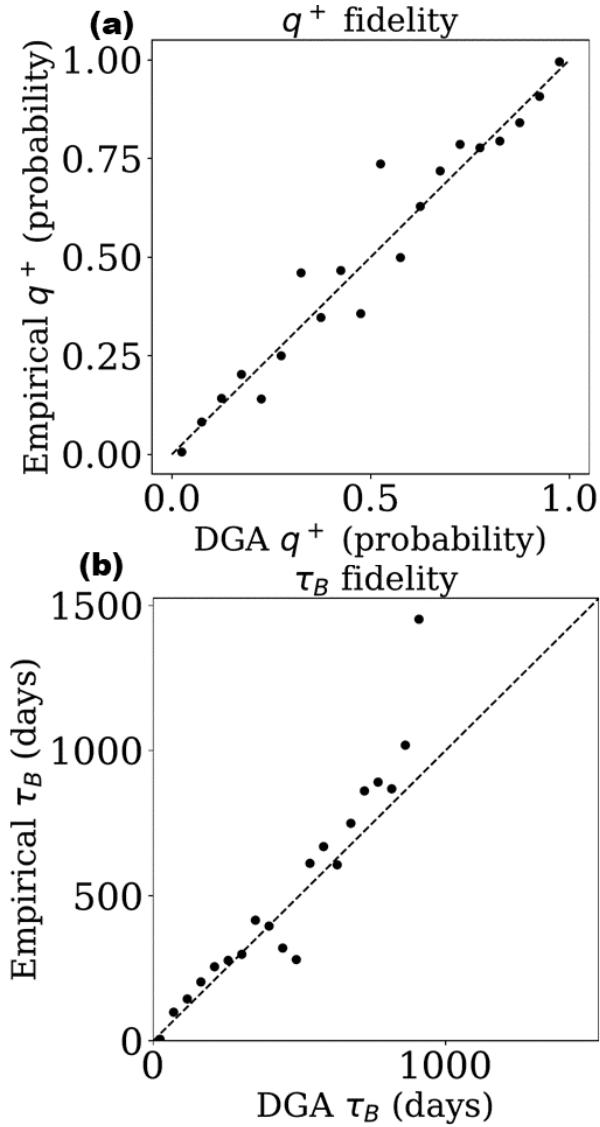
1097 **Fig. 7. Results of LASSO regression of the forward committor with  $U, \text{Re}\Psi, \text{Im}\Psi$  as input**  
 1098 **features.** Panel (a) shows the coefficients when  $q^+$  is regressed as a function of only the  
 1099 variables at a given altitude, and panel (b) shows the corresponding correlation score. 21.5  
 1100 km seems the most predictive (where  $z \equiv 0$  at the tropopause, not the surface). Panel (c)  
 1101 shows the coefficient structure when all altitudes are considered simultaneously. By design,  
 1102 most of the coefficients are zero, but most of the nonzero coefficients appear at 21.5 km, once  
 1103 again distinguishing that level as highly relevant for prediction. . . . . 60

1104 **Fig. 8. Two-dimensional projections of the committor and mean first passage times.** We have  
 1105 projected three quantities onto the observable subspace of zonal wind and imaginary part of  
 1106 streamfunction at 21.5 km. (a) forward committor  $q^+(x) = \mathbb{P}_x\{\tau_B < \tau_A\}$ , (b) first passage  
 1107 time to  $B$   $\mathbb{E}[\tau_B]$ , and (c) conditional mean first passage time to  $B$   $\mathbb{E}[\tau_B|\tau_B < \tau_A]$ . The  
 1108 condition  $\tau_B < \tau_A$  decreases the passage time by an order of magnitude, because it excludes  
 1109 the possibility of getting trapped in  $A$  first. Figure 9 quantifies the relationship between the  
 1110 committor and conditional passage time, and its forecasting implications. . . . . 61

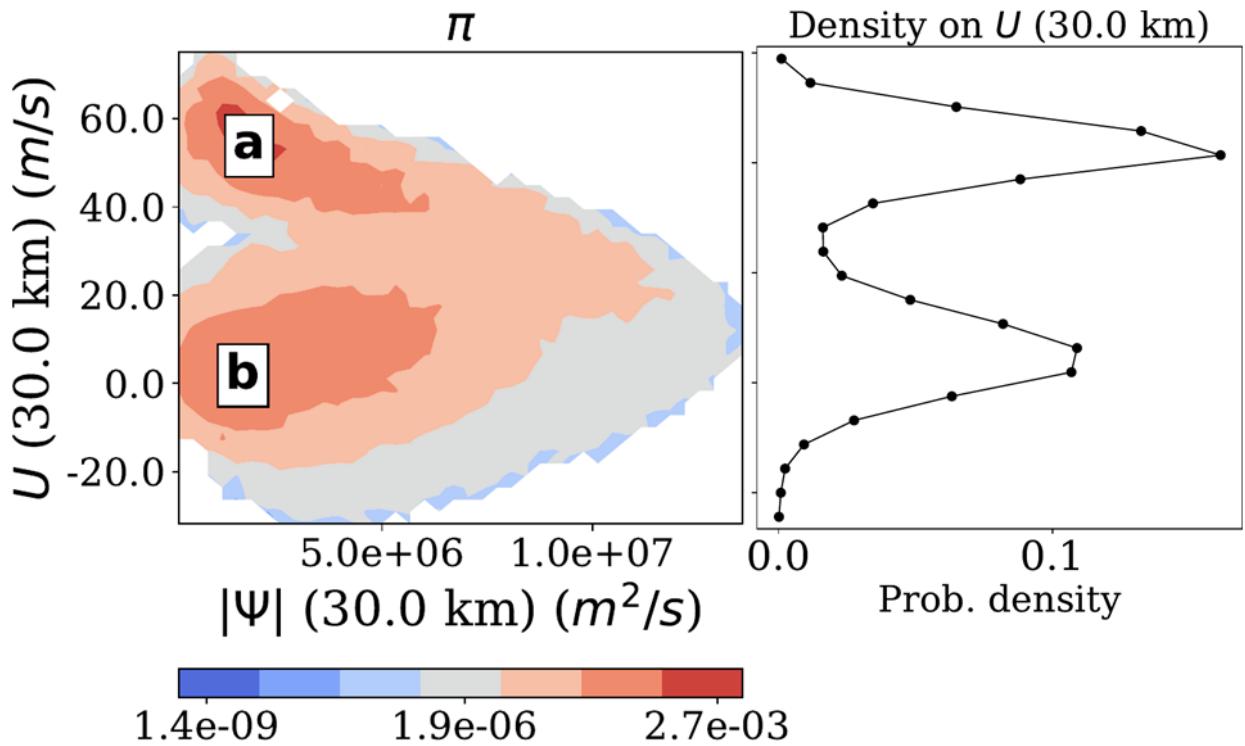
1111 **Fig. 9. Relationship between committor and mean first passage time.** Panel (a) shows the  
 1112 relationship between  $q^+$  (probability of next hitting  $B$ ) and  $\mathbb{E}[\tau_B|\tau_B < \tau_A]$ , the time until  
 1113 hitting  $B$  conditional on avoiding  $A$ . These quantities correspond to panels (a) and (c) of  
 1114 Figure 8. Panel (b) shows the same relationship but in the  $B \rightarrow A$  direction. In both cases, we  
 1115 performed a least squares regression weighted by the change of measure. A +0.1 increase  
 1116 in the probability  $q^+$  of next hitting  $B$  comes with a 6.3-day decrease in the expected time to  
 1117 get there, whereas a +0.1 increase in the opposite probability  $1 - q^+$  comes with a 9.8-day  
 1118 reduction in the time to reach  $A$ . Meanwhile, the vertical intercepts indicate the mean time  
 1119 of a full transition from  $A \rightarrow B$  (79 days) and  $B \rightarrow A$  (170 days). . . . . 62



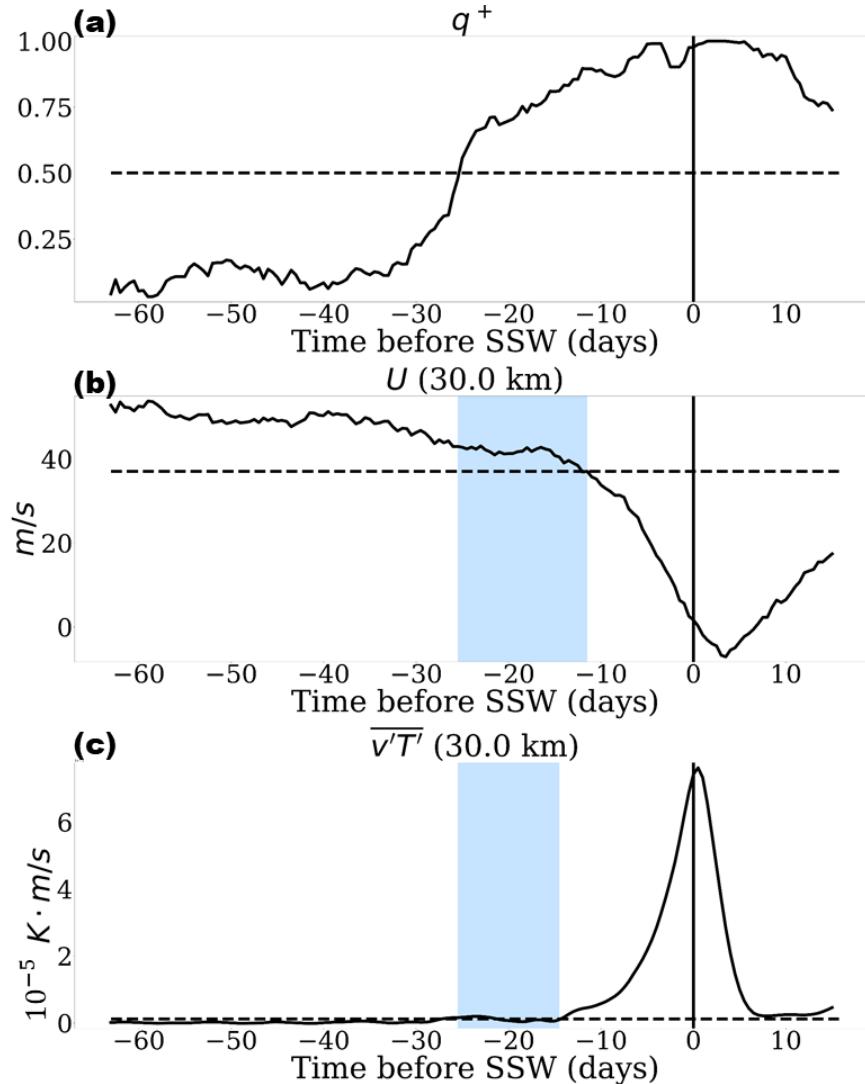
1120 FIG. 1. **Illustration of the two stable states of the Holton-Mass model and transitions between them.** (a)  
 1121 Zonal wind profiles of the radiatively maintained strong vortex (the fixed point **a**, blue) which increases linearly  
 1122 with altitude, and the weak vortex (the fixed point **b**, red) which dips close to zero in the mid-stratosphere. (b)  
 1123 Streamfunction contours are overlaid for the two equilibria **a** and **b**, the weak vortex exhibiting strong westward  
 1124 phase tilt with altitude. (c) Timeseries of  $U(30\text{ km})$  from a long stochastic simulation, including several noise-  
 1125 induced transitions from *A* to *B* (green) and from *B* to *A* (orange). Although both states **a** and **b** are equilibria in  
 1126 this parameter regime ( $h = 38.5m$ ), the stochastic perturbations uncover the vacillation cycles that would appear  
 1127 beyond the Hopf bifurcation if  $h$  were increased. (d) A parametric curve of the same trajectory segment as in  
 1128 (c) with the same color highlights for transition paths, but in the space  $(|\Psi|, U)$  at 30 km. The two equilibria are  
 1129 indicated with horizontal blue and red lines.



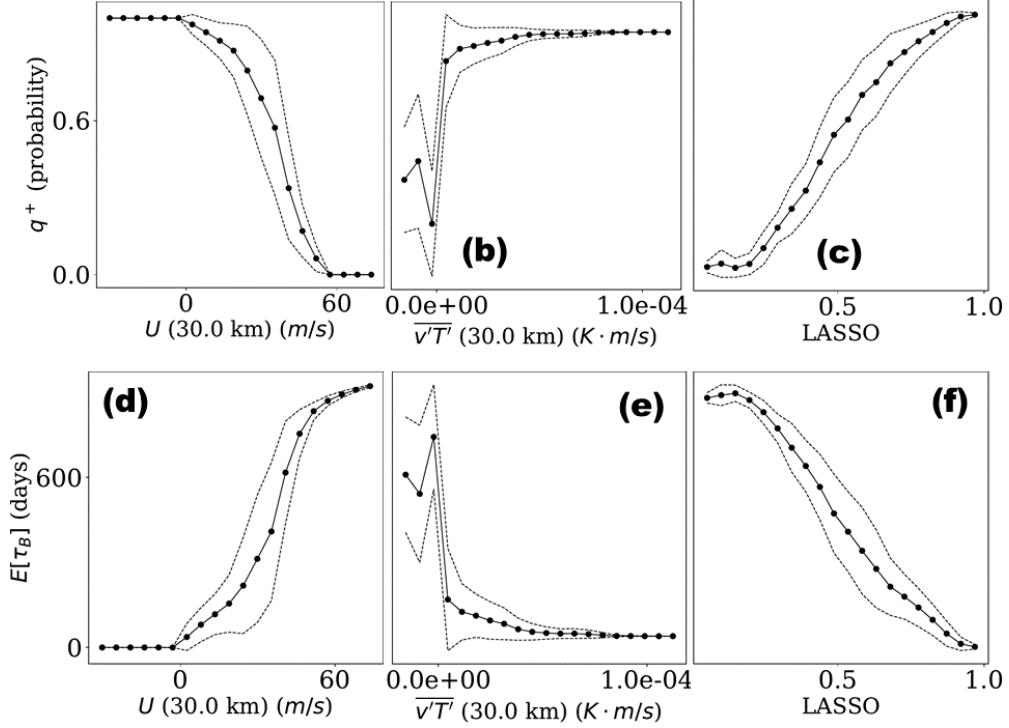
1130 FIG. 2. **Accuracy of the committor and mean first passage time calculations verified with long trajectory**  
 1131 **data.** The DGA calculations assign approximate committor and mean first passage time values  $q^+(\mathbf{X}_n(0))$  and  
 1132  $m_B(\mathbf{X}_n(0))$  to each data point. Because each snapshot  $x_n$  was collected from a long trajectory, its destination and  
 1133 the time to get there are known and can provide an empirical validation of the committor and lead time. For 20  
 1134 equal partitions  $(\zeta_1, \zeta_2)$  of the interval  $(0, 1)$ , we assemble all trajectory starts  $\mathbf{X}_n(0)$  with  $q^+(\mathbf{X}_n(0)) \in (\zeta_1, \zeta_2)$   
 1135 and count the fraction heading toward  $B$ . These are the empirical committors for the interval  $(\zeta_1, \zeta_2)$ , and are  
 1136 plotted on the vertical axis against  $(\zeta_1 + \zeta_2)/2$ . Similarly, we bin the space of calculated first passage times and for  
 1137 each bin average the empirical first passage time to  $B$ . Both quantities line up well between DGA computations  
 1138 and empirical values, with the exception of the longest passage times, which are underestimated.



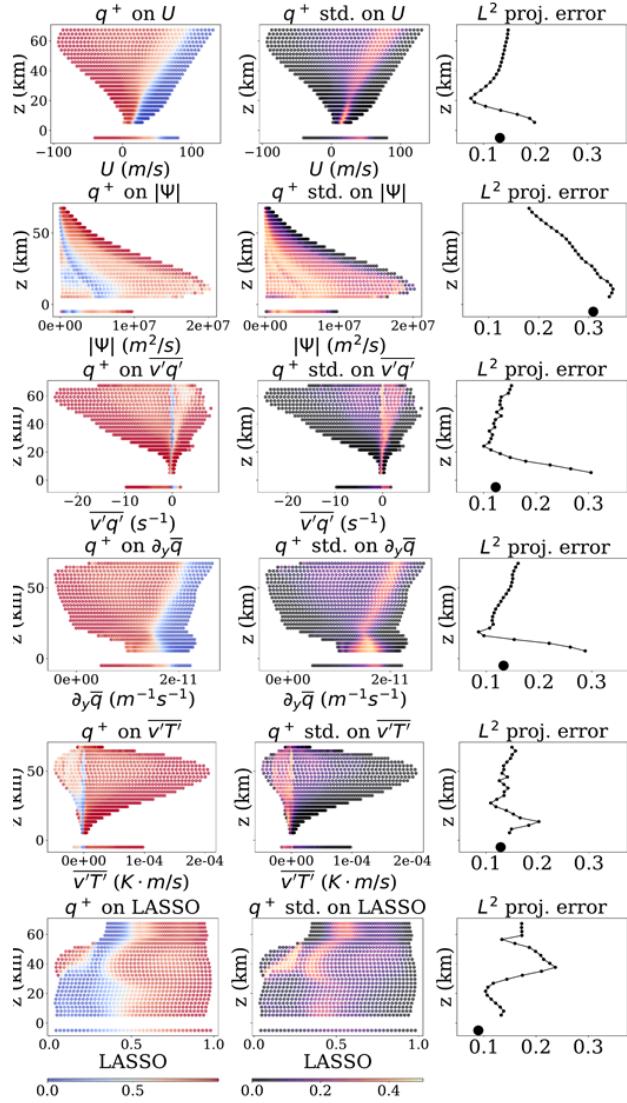
1139 FIG. 3. **Steady-state distribution.** The density  $\pi(\mathbf{x})$  is projected onto the two-dimensional space  $(|\Psi|, U)$  at  
 1140 30 km, on a log scale. The density is peaked in the neighborhoods of the two fixed points. On the right is a  
 1141 projection of  $\pi$  onto the single variable  $U(30 \text{ km})$ , on a linear scale, confirming strong bimodality.



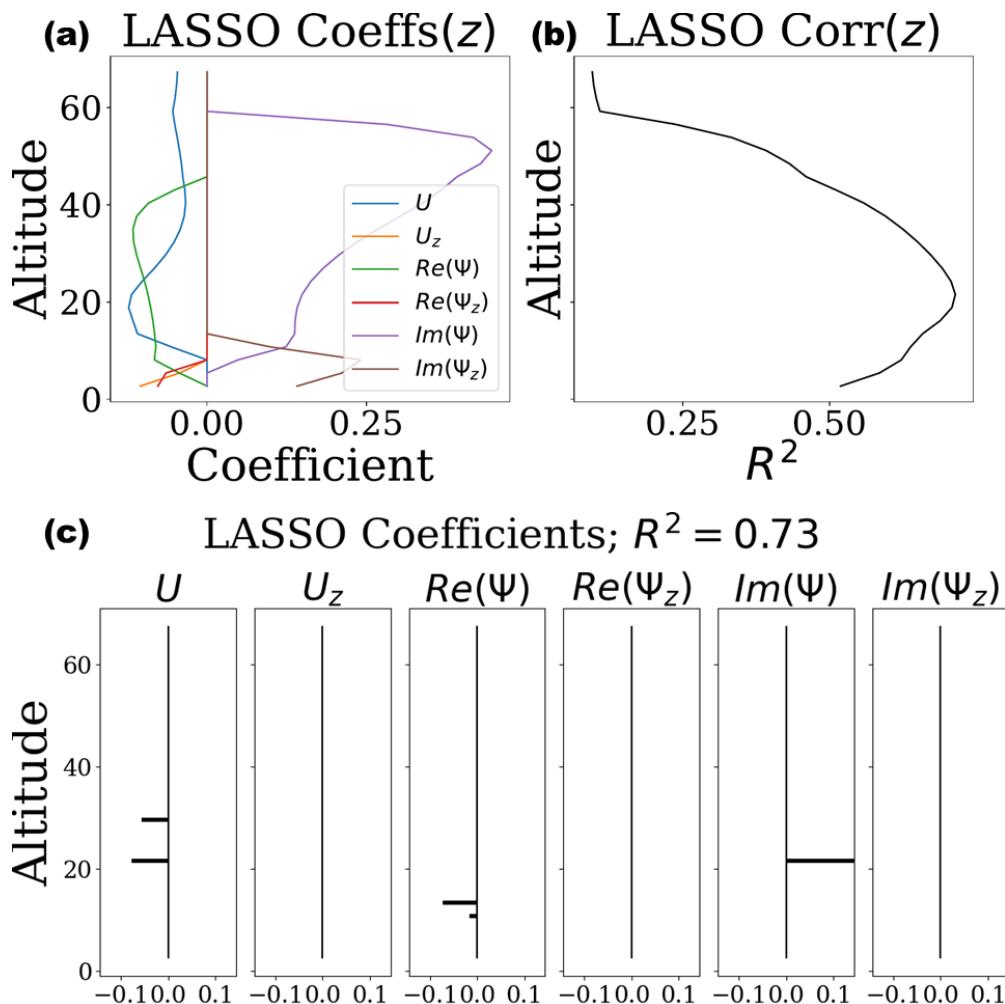
1142 **FIG. 4. The committor vs. other observables as a forecasting tool.** A representative simulated SSW event  
 1143 from the long simulation is plotted over time, starting 65 days in advance of the official event when  $U(30km)$   
 1144 first drops below 1.75 m/s, which is marked by a vertical solid line. Panel (a) shows the committor over time  
 1145 following the trajectory, panel (b) shows the zonal wind  $U(30km)$ , and panel (c) shows the eddy heat flux  
 1146  $\overline{v'T'}$ (30 km). Horizontal dashed lines mark the natural forecasting threshold of  $q^+ = 0.5$  (panel (a)) or the value  
 1147 of the observable most closely associated with  $q^+ = 0.5$ : 37 m/s (panel (b)) and  $1.2 \times 10^{-6} K \cdot m/s$  (panel (c)).  
 1148 The sharp increase in  $q^+$  as it crosses the threshold provides a clear and early warning sign of oncoming SSW,  
 1149 about 26 days in advance.  $U$  and  $\overline{v'T'}$  are moving slowly at that time, and don't clear their respective thresholds  
 1150 for the last time until the event is much closer at hand. The gap in lead time is marked by blue strips.



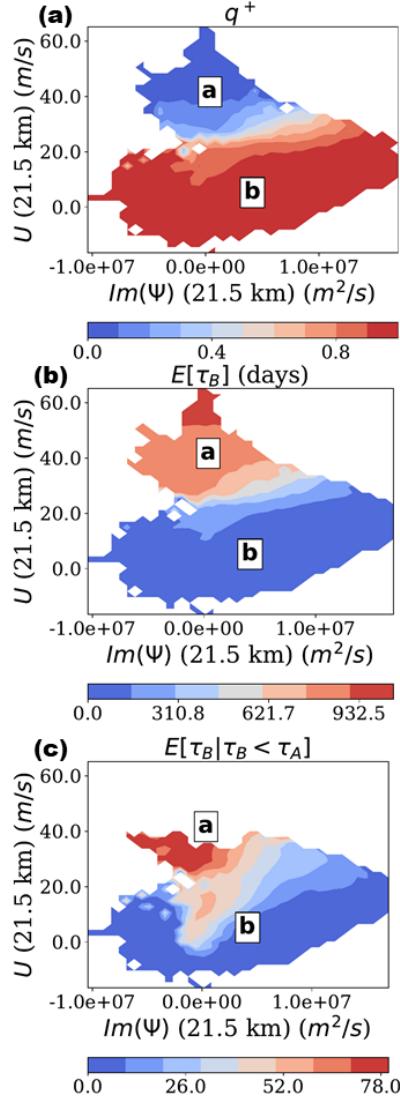
1151 **FIG. 5. One-dimensional projections of the forward committor and mean first passage time to  $B$ ,**  
 1152 **computed with DGA.** These functions depend on all  $d = 75$  degrees of freedom in the model, but we have  
 1153 averaged across  $d - 1 = 74$  dimensions to visualize the committor (first row) and mean first passage time to  
 1154  $B$  (second row) as rough functions of three single degrees of freedom:  $U(30\text{ km})$  (first column),  $\overline{v'T'}$  (30 km)  
 1155 (second column), and the LASSO-regressed committor (third column). The forward committor measures  
 1156 proximity to  $B$  in probability, while mean passage time to  $B$  measures proximity in time, hence the negative  
 1157 correlation between the two quantities. The general trends reveal fairly obvious relationships: stronger wind  
 1158 is associated with tendency towards the strong-vortex state  $A$ , and larger poleward eddy heat flux is associated  
 1159 with tendency toward the weak vortex state  $B$ . In addition, curves like this assess the quality of single-variable  
 1160 observables as proxies for an oncoming transition event. The committor and passage time vary smoothly and  
 1161 (mostly) monotonically with  $U$ , but discontinuously with  $\overline{v'T'}$ : the heat flux burst that accompanies a SSW gives  
 1162 no advance warning for the event, while a small negative change in  $U$  indicates incrementally higher transition  
 1163 probability and shorter lead time.



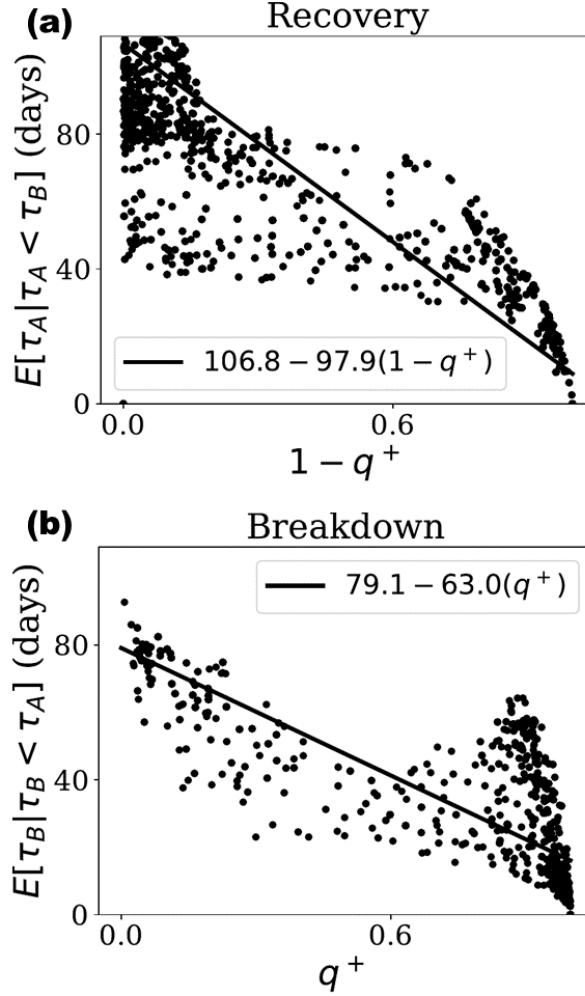
1164 **FIG. 6. Projection of the forward committor onto a large collection of one-dimensional CVs, along with the**  
 1165 **associated standard deviation, or projection error, of the committor along the remaining 74 model dimensions.**  
 1166 Consider the first two panels. The left-hand panel shows, for each discretized altitude  $z$ , a heatmap of the  
 1167 committor as  $U(z)$  ranges from its minimum to its maximum realized strength at that altitude. At the bottom  
 1168 is an additional heatmap of the committor as a function of  $\langle U \rangle_z$ , the vertical average. These are conditional  
 1169 expectations, with the corresponding conditional standard deviations displayed in the right-hand panels. The  
 1170 following rows display analogous plots for wave magnitude, eddy PV flux, background PV gradient, eddy heat  
 1171 flux, and the LASSO predictor specified in Figure 7. The bottom line on the last plot is not a vertical average,  
 1172 but the results of regression on all altitudes at once.



1173 FIG. 7. Results of LASSO regression of the forward committor with  $U$ ,  $Re\Psi$ ,  $Im\Psi$  as input features. Panel  
 1174 (a) shows the coefficients when  $q^+$  is regressed as a function of only the variables at a given altitude, and panel  
 1175 (b) shows the corresponding correlation score. 21.5 km seems the most predictive (where  $z \equiv 0$  at the tropopause,  
 1176 not the surface). Panel (c) shows the coefficient structure when all altitudes are considered simultaneously. By  
 1177 design, most of the coefficients are zero, but most of the nonzero coefficients appear at 21.5 km, once again  
 1178 distinguishing that level as highly relevant for prediction.



1179 **FIG. 8. Two-dimensional projections of the committor and mean first passage times.** We have projected  
 1180 three quantities onto the observable subspace of zonal wind and imaginary part of streamfunction at 21.5 km.  
 1181 (a) forward committor  $q^+(x) = \mathbb{P}_x\{\tau_B < \tau_A\}$ , (b) first passage time to  $B$   $\mathbb{E}[\tau_B]$ , and (c) conditional mean first  
 1182 passage time to  $B$   $\mathbb{E}[\tau_B | \tau_B < \tau_A]$ . The condition  $\tau_B < \tau_A$  decreases the passage time by an order of magnitude,  
 1183 because it excludes the possibility of getting trapped in  $A$  first. Figure 9 quantifies the relationship between the  
 1184 committor and conditional passage time, and its forecasting implications.



1185 FIG. 9. **Relationship between committor and mean first passage time.** Panel (a) shows the relationship  
 1186 between  $q^+$  (probability of next hitting  $B$ ) and  $\mathbb{E}[\tau_B | \tau_B < \tau_A]$ , the time until hitting  $B$  conditional on avoiding  
 1187  $A$ . These quantities correspond to panels (a) and (c) of Figure 8. Panel (b) shows the same relationship but in the  
 1188  $B \rightarrow A$  direction. In both cases, we performed a least squares regression weighted by the change of measure. A  
 1189  $+0.1$  increase in the probability  $q^+$  of next hitting  $B$  comes with a 6.3-day decrease in the expected time to get  
 1190 there, whereas a  $+0.1$  increase in the opposite probability  $1 - q^+$  comes with a 9.8-day reduction in the time to  
 1191 reach  $A$ . Meanwhile, the vertical intercepts indicate the mean time of a full transition from  $A \rightarrow B$  (79 days) and  
 1192  $B \rightarrow A$  (170 days).