

1 **Predicting Observed and Hidden Extreme Events in Complex Nonlinear Dynamical**
2 **Systems with Partial Observations and Short Training Time Series**

3 Nan Chen^{1, a)} and Andrew J. Majda^{2, b)}

4 ¹⁾*Department of Mathematics, University of Wisconsin-Madison*

5 ²⁾*Department of Mathematics and Center for Atmosphere Ocean Science,*
6 *Courant Institute of Mathematical Sciences, New York University, New York, NY,*
7 *USA, and*

8 *Center for Prototype Climate Modeling, New York University Abu Dhabi,*
9 *Saadiyat Island, Abu Dhabi, UAE*

10 (Dated: 30 July 2019)

11 Extreme events appear in many complex nonlinear dynamical systems. Predicting
12 extreme events has important scientific significance and large societal impacts. In
13 this paper, a new mathematical framework of building suitable nonlinear approxi-
14 mate models is developed, which aims at predicting both the observed and hidden
15 extreme events in complex nonlinear dynamical systems for short-, medium- and
16 long-range forecasting using only short and partially observed training time series.
17 Different from many ad-hoc data-driven regression models, these new nonlinear mod-
18 els take into account physically motivated processes and physics constraints. They
19 also allow efficient and accurate algorithms for parameter estimation, data assimi-
20 lation and prediction. Cheap stochastic parameterizations, judicious linear feedback
21 control and suitable noise inflation strategies are incorporated into the new nonlin-
22 ear modeling framework, which provide accurate predictions of both the observed
23 and hidden extreme events as well as the strongly non-Gaussian statistics in various
24 highly intermittent nonlinear dyad and triad models, including the Lorenz 63 model.
25 Then a stochastic mode reduction strategy is applied to a 21-dimensional nonlinear
26 paradigm model for topographic mean flow interaction. The resulting 5-dimensional
27 physics-constrained nonlinear approximate model is able to accurately predict the ex-
28 treme events and the regime switching between zonally blocked and unblocked flow
29 patterns. Finally, incorporating judicious linear stochastic processes into a simple
30 nonlinear approximate model succeeds in learning certain complicated nonlinear ef-
31 fects of a 6-dimensional low-order Charney-DeVore model with strong chaotic and
32 regime switching behavior. The simple nonlinear approximate model then allows
33 accurate online state estimation and the short- and medium-range forecasting of ex-
34 treme events.

35 PACS numbers: 02.50.Ey, 05.45.-a, 05.45.Tp, 92.60.Aa

a) chenman@math.wisc.edu; Corresponding author

b) jonjon@cims.nyu.edu

36 Extreme events appear in many complex nonlinear dynamical systems. These
37 extreme events are associated with the sudden changes of states in the underlying
38 complex systems and the occurrence of extreme events often results in large
39 social impact. Therefore, predicting extreme events has both scientific significance
40 and practical implications. However, the big challenges of predicting the
41 extreme events in complex nonlinear systems include the lack of understanding
42 of physics, the huge computational cost in running the complex models
43 and data assimilation, as well as the availability of only short and partially observed
44 training data. In this paper, a new mathematical framework of building
45 suitable nonlinear approximate models is developed, which aims at predicting
46 both the observed and hidden extreme events in complex nonlinear dynamical
47 systems for short-, medium- and long-range forecasting using only short and
48 partially observed training time series. This framework also allows efficient and
49 accurate data assimilation, parameter estimation and prediction algorithms. Different
50 effective and practical strategies are incorporated into the framework to
51 develop suitable approximate models for predicting extreme events and other
52 non-Gaussian features in various complex turbulent dynamical systems.

53 I. INTRODUCTION

54 Extreme events appear in many complex nonlinear dynamical systems in geoscience, engineering,
55 excitable media, neural science and material science¹⁻⁸. Examples include oceanic
56 rogue waves^{9,12}, extreme weather and climate patterns^{10,11} such as blocking events and
57 turbulent tracers¹³⁻¹⁵, and bursting neurons¹⁶. These extreme events are associated with the
58 sudden changes of states in the underlying complex systems and the occurrence of extreme
59 events often results in large social impact. Therefore, predicting extreme events has both
60 scientific significance and practical implications.

61 However, predicting the extreme events in complex nonlinear systems is quite challenging.
62 First, nature or the perfect model is never known in practice. Model error due to the lack of
63 the understanding of physics may prevent the skillful predictions of the extreme events^{1,17-19}.
64 Second, even if the perfect model is known, the underlying nonlinear dynamics of nature
65 can be extremely complicated with strong non-Gaussian characteristics, multiscale features

66 and high dimensionality^{10,20,21}. Thus, running the perfect model is usually computationally
67 unaffordable for real-time prediction. On the other hand, despite that coarse-graining the
68 numerical resolutions improves the computational efficiency, such a numerical approximation
69 often results in missing the key nonlinear interactions between different temporal and spatial
70 scales and brings about large errors, especially for extreme events. Third, it is important to
71 notice that only partial and noisy observations are available in many practical situations^{22–24},
72 which implies the states of the unobserved variables have to be estimated via online data
73 assimilation algorithms. Unfortunately, the existing data assimilation algorithms for general
74 complex nonlinear dynamical systems are either quite expensive (e.g., particle filter) or
75 involving intrinsic approximate errors due to the coarse-grained statistics (e.g., ensemble
76 Kalman filter)^{21,25–28}. The assimilated states from the latter may also contain large biases
77 due to the fact that high order moments are important contributors to the extreme events.
78 Finally, the actual climate signal is often measured through time series. However, since the
79 high-resolution satellites and other refined measurements were not widely developed until
80 recent times, the available useful training data is very limited with about only 50 years in
81 many real applications. Thus, predicting extreme events using short and partially observed
82 training time series is another remarkably challenging task.

83 For the reasons given above, developing suitable approximate models for predicting ex-
84 treme events is crucial in practice. These approximate models aim at capturing the key
85 nonlinear dynamical and non-Gaussian statistical features of nature. They also need to be
86 computationally tractable and allow efficient algorithms for online data assimilation, pa-
87 rameter estimation and prediction. There have been some recent progress in the extreme
88 events prediction. For example, a new statistical dynamical model was developed to predict
89 extreme events and anomalous features in shallow water waves¹². A suite of reduced-order
90 stochastic models was built, which succeeds in predicting the extreme events in complex
91 geophysical flows²⁹ and their long-term non-Gaussian features³⁰ as well as forecasting the
92 associated statistical responses and quantifying the uncertainty³¹. In addition, mode decom-
93 position techniques were applied for probing the most unstable modes and building low-order
94 models for extreme events prediction^{32,33}.

95 In this paper, a new mathematical framework of building suitable nonlinear approximate
96 models is developed, which aims at predicting both the observed and hidden extreme events
97 in complex nonlinear dynamical systems using only short and partially observed training

98 time series. The models belonging to this mathematical framework are highly nonlinear and
99 are able to capture many key non-Gaussian characteristics as observed in nature³⁴. Unlike
100 traditional regression and other ad hoc models with prescribed basis functions or structures,
101 this framework contains a rich class of statistical dynamical models and is amenable to a
102 wide range of applications. One important feature of this nonlinear modeling framework is
103 that physically motivated processes and physics constraints^{35,36} can be incorporated into the
104 models, which is fundamentally different from many purely data-driven statistical models
105 that have no clear physical meanings. Such a trait not only enables the models to take
106 into account both the dynamical and statistical information but also allows using only a
107 short training time series for model calibration. The latter is due to the (partially) iden-
108 tified dynamical structures from some physics reasoning and physics constraints. Another
109 key advantage of this new framework is that despite the intrinsic nonlinearity, it allows
110 closed analytic formulae for assimilating the states of the unobserved variables^{37,38}, which is
111 computationally efficient and accurate. This provides an extremely useful and practical ap-
112 proach for predicting extreme events and other non-Gaussian features in complex nonlinear
113 dynamical systems.

114 Short-, medium- and long-range forecasting of extreme events all have practical significance^{10,39-41}.
115 The efficient data assimilation scheme associated with the nonlinear models within the above
116 framework provides an accurate estimation of the initial values, which play a crucial role in
117 improving the short-term prediction skill. On the other hand, the focus of the long-term
118 prediction is on the statistics, which is calculated by making use of a long trajectory together
119 with the ergodic property of many complex turbulent systems¹. In particular, reproducing
120 the statistical equilibrium non-Gaussian probability density function (PDF) with fat tails is
121 a good evidence of the successful prediction of extreme events, where the extreme events and
122 intermittency are the main contributors to the fat tails. The medium-range forecast aims at
123 recovering the transition behavior of the underlying dynamics. A skillful medium-range pre-
124 diction requires both an accurate estimation of the initial values and a suitable description
125 of the time evolution of the approximate model, and is often a challenging task. Finally,
126 certain internal or external perturbations are able to kick the model variables outside the
127 attractor. Therefore, predicting the time evolution of the extreme events that start from a
128 state outside the attractor also has practical importance. It is worth remarking that many
129 purely data-driven or machine learning methods fail to predict extreme events even though

130 most of those methods show high skill in fitting the observed time series. For example,
131 as has been pointed out in some recent work^{42,43}, even one of the most advanced neural
132 networks with long short-term memory⁴⁴ and the Gaussian process regression⁴⁵ suffer from
133 a finite time blowup issue when they are applied for predicting extreme events. Such a
134 pathological behavior can only be overcome by using hybrid strategies that combine these
135 methods with suitable models⁴². Note that these purely data-driven methods often demand
136 tremendous training data^{46,47}, which is not practical in many scientific scenarios where only
137 short training time series are available. In addition, without suitable models, predicting
138 extreme events in the unobserved processes becomes extremely difficult.

139 This paper aims at incorporating practical strategies into the development of suitable
140 approximate models for predicting both the observed and hidden extreme events. These
141 approximate models belong to the new nonlinear modeling framework, which allows an ef-
142 ficient and accurate data assimilation scheme and only short and partially observed time
143 series are needed for model calibration. The first effective strategy is to adopt simple s-
144 tochastic parameterizations for approximating complicated hidden processes. Despite the
145 simple forms, the judicious applications of these stochastic parameterizations are neverthe-
146 less able to capture the nonlinear interactions between the observed and hidden variables
147 and predict the associated extreme events. Such an idea has been successfully applied to
148 the stochastic parameterized extended Kalman filter (SPEKF) forecast models^{48,49}, dynamic
149 stochastic superresolution of sparsely observed turbulent systems^{50,51} and stochastic super-
150 parameterization for geophysical turbulence⁵². The second strategy here is motivated from
151 control theory, which involves incorporating simple feedback control terms into the approx-
152 imate models for model simplification. This simple feedback control strategy succeeds in
153 capturing the key nonlinear statistical interactions as well as the causal effects between
154 the observed and hidden variables, which are essential to accurately predicting the extreme
155 events in the hidden processes. Note that predicting the hidden extreme events is typical-
156 ly a great challenge given only partial observations. The third strategy makes use of the
157 stochastic mode reduction technique⁵³⁻⁵⁶, which allows a significant dimension reduction in
158 the approximate models for many multiscale turbulent dynamical systems while the reduced
159 order models retain the crucial nonlinear and non-Gaussian features. Applying this strategy,
160 the nonlinear effects of the unresolved fast modes in the motion of the resolved variables
161 are represented by effective damping and stochastic forcing. The resulting approximate

162 models naturally belong to the new nonlinear modeling framework that allows extremely
163 efficient data assimilation and prediction schemes. These approximate models also preserve
164 physics-constrained properties. Another extremely useful strategy is to incorporate simple
165 stochastic processes with additive noise and memory into the approximate models, which
166 aim at effectively describing certain complicated nonlinear components that are hard to
167 deal with in strongly nonlinear and chaotic dynamical systems. Due to the unique feature
168 of the new nonlinear modeling framework, it allows an efficient and accurate way of using
169 simple stochastic processes to learn these complex nonlinear components on the fly, which
170 greatly facilitates the short- and medium-range forecasts of both the observed and hidden
171 extreme events. Other approaches of building approximate predictive models that can be
172 incorporated into the new nonlinear framework developed here involve using the noise in-
173 flation technique to effectively characterize the contributions from small-scale variables and
174 fast-wave averaging of the variables with rapid decaying⁵⁷.

175 The rest of the paper is organized as follows. Section II describes the new nonlinear
176 mathematical framework for developing suitable approximate models. Section III contains
177 the efficient and accurate data assimilation, parameter estimation and prediction algorithms.
178 Both the path-wise and information measurements in quantifying the prediction skill are also
179 included in this section. Section IV illustrates the skill of predicting intermittent extreme
180 events using cheap stochastic parameterizations with significant model error. Section V
181 makes use of a nonlinear energy-conserving dyad model to show the success of applying
182 the simple feedback control strategy in facilitating the prediction of the hidden extreme
183 events. The effect of noise inflation in approximate models for predicting extreme events is
184 illustrated based on the chaotic Lorenz 63 model in Section VI. Section VII starts with a
185 21-dimensional nonlinear topographic mean flow interaction model with regime switching.
186 Stochastic mode reduction strategy is applied in a suitable way to develop an approximate
187 nonlinear model with only 5 dimensions, which is nevertheless able to predict the observed
188 and hidden extreme events as well as the regime switching between zonally blocked and
189 unblocked flow patterns with high accuracy. In Section VIII, it is shown that incorporating
190 judicious linear stochastic processes into a simple nonlinear approximate model succeeds in
191 learning certain complicated nonlinear effects of a 6-dimensional low-order chaotic Charney-
192 DeVore model with strong chaotic and regime switching behavior. The resulting nonlinear
193 approximate model allows accurate online state estimation and the short- and medium-range

194 forecasting of extreme events. The paper is concluded in Section IX.

195 II. A NONLINEAR MATHEMATICAL MODELING FRAMEWORK 196 WITH SOLVABLE CONDITIONAL STATISTICS

A nonlinear mathematical modeling framework is established in this section, which will be used to the development of suitable approximate models for predicting extreme events. The general form of the nonlinear models within this framework is the following³⁸,

$$d\mathbf{u}_I = [\mathbf{A}_0(t, \mathbf{u}_I) + \mathbf{A}_1(t, \mathbf{u}_I)\mathbf{u}_{II}]dt + \Sigma_I(t, \mathbf{u}_I)d\mathbf{W}_I(t), \quad (1a)$$

$$d\mathbf{u}_{II} = [\mathbf{a}_0(t, \mathbf{u}_I) + \mathbf{a}_1(t, \mathbf{u}_I)\mathbf{u}_{II}]dt + \Sigma_{II}(t, \mathbf{u}_I)d\mathbf{W}_{II}(t), \quad (1b)$$

197 where the state variables are written in the form $\mathbf{u} = (\mathbf{u}_I, \mathbf{u}_{II})$ with both $\mathbf{u}_I \in R^{N_I}$ and
198 $\mathbf{u}_{II} \in R^{N_{II}}$ being multidimensional variables. In (1), $\mathbf{A}_0, \mathbf{A}_1, \mathbf{a}_0, \mathbf{a}_1, \Sigma_I$ and Σ_{II} are vectors
199 and matrices that depend only on time t and the state variables \mathbf{u}_I , and $\mathbf{W}_I(t)$ and $\mathbf{W}_{II}(t)$
200 are independent Wiener processes. The systems in (1) are named as conditional Gaussian
201 systems due to the fact that once $\mathbf{u}_I(s)$ for $s \leq t$ is given, $\mathbf{u}_{II}(t)$ conditioned on $\mathbf{u}_I(s)$
202 becomes a Gaussian process with mean $\bar{\mathbf{u}}_{II}(t)$ and covariance $\mathbf{R}_{II}(t)$, i.e.,

$$p(\mathbf{u}_{II}(t)|\mathbf{u}_I(s \leq t)) \sim \mathcal{N}(\bar{\mathbf{u}}_{II}(t), \mathbf{R}_{II}(t)). \quad (2)$$

203 Despite the conditional Gaussianity, the coupled system (1) remains highly nonlinear and
204 is able to capture the non-Gaussian features as in nature. This conditional Gaussian nonlin-
205 ear modeling framework includes many physics-constrained nonlinear stochastic models^{35,36},
206 large-scale dynamical models in turbulence, fluids and geophysical flows, as well as stochas-
207 tically coupled reaction-diffusion models in neuroscience and ecology. See a recent work³⁴
208 for a gallery of examples of the conditional Gaussian systems. Applications of the condition-
209 al Gaussian systems to strongly nonlinear systems include developing low-order nonlinear
210 stochastic models for predicting the intermittent time series of the Madden-Julian oscillation
211 (MJO) and the monsoon intraseasonal variabilities⁵⁸⁻⁶⁰, filtering the stochastic skeleton
212 model for the MJO⁶¹, and recovering the turbulent ocean flows with noisy observations from
213 Lagrangian tracers⁶²⁻⁶⁴. Other studies that also fit into the conditional Gaussian framework
214 includes the cheap exactly solvable forecast models in dynamic stochastic superresolution of
215 sparsely observed turbulent systems^{50,51}, stochastic superparameterization for geophysical
216 turbulence⁵² and blended particle filters for large-dimensional chaotic systems⁶⁵.

One important feature of the above conditional Gaussian nonlinear framework is that the conditional Gaussian distribution $p(\mathbf{u}_{\text{II}}(t)|\mathbf{u}_{\text{I}}(s \leq t))$ in (2) has closed analytic form³⁷,

$$d\bar{\mathbf{u}}_{\text{II}}(t) = [\mathbf{a}_0(t, \mathbf{u}_{\text{I}}) + \mathbf{a}_1(t, \mathbf{u}_{\text{I}})\bar{\mathbf{u}}_{\text{II}}]dt + (\mathbf{R}_{\text{II}}\mathbf{A}_1^*(t, \mathbf{u}_{\text{I}}))(\boldsymbol{\Sigma}_{\text{I}}\boldsymbol{\Sigma}_{\text{I}}^*)^{-1}(t, \mathbf{u}_{\text{I}}) \times [d\mathbf{u}_{\text{I}} - (\mathbf{A}_0(t, \mathbf{u}_{\text{I}}) + \mathbf{A}_1(t, \mathbf{u}_{\text{I}})\bar{\mathbf{u}}_{\text{II}})dt], \quad (3a)$$

$$d\mathbf{R}_{\text{II}}(t) = \left\{ \mathbf{a}_1(t, \mathbf{u}_{\text{I}})\mathbf{R}_{\text{II}} + \mathbf{R}_{\text{II}}\mathbf{a}_1^*(t, \mathbf{u}_{\text{I}}) + (\boldsymbol{\Sigma}_{\text{II}}\boldsymbol{\Sigma}_{\text{II}}^*)(t, \mathbf{u}_{\text{I}}) - (\mathbf{R}_{\text{II}}\mathbf{A}_1^*(t, \mathbf{u}_{\text{I}}))(\boldsymbol{\Sigma}_{\text{I}}\boldsymbol{\Sigma}_{\text{I}}^*)^{-1}(t, \mathbf{u}_{\text{I}})(\mathbf{R}_{\text{II}}\mathbf{A}_1^*(t, \mathbf{u}_{\text{I}}))^* \right\} dt. \quad (3b)$$

217 It is natural to assume \mathbf{u}_{I} contains the observed variables while \mathbf{u}_{II} is a collection of the
 218 unobserved ones. Therefore, the analytically solvable conditional statistics in (3) allows an
 219 extremely efficient and accurate way of estimating the hidden states given the observations,
 220 known as the data assimilation, which facilitates predictions. Note that in the data assimi-
 221 lation language the conditional mean and conditional covariance in (3) are also known as the
 222 posterior mean and posterior covariance. In addition, the conditional Gaussian nonlinear
 223 modeling framework (1) and its closed analytical form of the conditional statistics (3) offer a
 224 statistical efficient and accurate way of solving the time evolution of the associated Fokker-
 225 Planck equation in high dimensions^{66–68}, which also provides a powerful tool for carrying
 226 out ensemble forecasts.

227 III. DATA ASSIMILATION, PREDICTION, AND THE 228 QUANTIFICATION OF PREDICTION SKILL

229 A. Data assimilation of the unobserved variables

230 Data assimilation (also known as state estimation or filtering)^{21,25–28}, a procedure of es-
 231 timating the states of the unobserved variables, is the precondition of predicting complex
 232 dynamical systems. In fact, data assimilation of the unobserved variables can also be re-
 233 garded as the online “prediction” of these variables due to the fact that the recovered states
 234 of the unobserved variables are given by combining the information in the dynamics with
 235 the values of the observed variables.

236 Data assimilation of the unobserved variables plays an important role in short- and
 237 medium-range forecasts. This is because the ensemble prediction algorithm requires running
 238 the model forward with the given initial values for all the state variables. Since there is no

239 direct observations of the hidden or unresolved variables, assimilating their initial states
 240 becomes a necessity part of the ensemble forecast. In practice, the data assimilation is often
 241 required in an “online” form in the sense that the states of the unobserved variables need
 242 to be estimated at each time instant as time evolves. Therefore, developing an efficient and
 243 accurate data assimilation method is a crucial first step for predicting nonlinear complex
 244 dynamical systems and the associated extreme events. However, the classical Kalman filter
 245 or its continuous form Kalman-Bucy filter^{69–71} works only for linear models. On the other
 246 hand, for assimilating general complex nonlinear dynamical systems, the particle filter is
 247 quite expensive and contains sampling error while the ensemble Kalman filter takes into
 248 account only the first two moments which may end up with large biases for assimilating
 249 extreme events.

250 The conditional Gaussian nonlinear modeling framework in Section II provides an ef-
 251 ficient way of estimating the states of the highly non-Gaussian hidden variables $\mathbf{u}_{\text{II}}(t)$ in
 252 the complex nonlinear dynamical systems given the observations up to the current time
 253 $\mathbf{u}_{\text{I}}(s \leq t)$. The closed analytic formula in (3) avoids numerical and sampling errors, and it
 254 results in an extremely efficient and accurate way of computing the optimal states of $\mathbf{u}_{\text{II}}(t)$.

255 B. Short-, medium- and long-range forecasting

256 Prediction problems have been described by Lorenz as falling into two categories^{72,73}.
 257 Problems that depend on the initial condition, such as short- to medium-range weather
 258 forecasting, are described as “predictions of the first kind”, while problems for predicting
 259 the longer-term climatology, are referred to as “predictions of the second kind”.

260 For short- and medium-range forecasts, the system starts from an initial time t_0 , where
 261 the initial values of the unobserved variables are determined by data assimilation. Then
 262 an ensemble prediction algorithm is applied by running the model forward up to a given
 263 time t_1 . Typically, t_1 is not quite far from t_0 and therefore the system has not completely
 264 lost its memory of the initial values. Therefore, a good state estimation of the initial values
 265 via data assimilation plays an important role in providing an accurate short- and medium-
 266 range forecasting skill. The ensemble mean, which is the average value of all the ensemble
 267 members, is often used as a predictor for the evolution of the trajectories and the ensemble
 268 spread measures the uncertainty in the ensemble mean forecasts. The difference between

269 short- and medium-range forecasts is that the prediction skill at very short lead time largely
270 depends on the accuracy of the initial values while both the dynamical structures and the
271 initial values will be essential in predicting the model transition behavior in the medium
272 range. Capturing the time evolution of the large bursts in intermittent time series with
273 small uncertainty is the goal of short- and medium-range forecasts of the extreme events.

274 On the other hand, for the long-range forecast where t_1 is much larger than t_0 , the system
275 will lose its memory from the initial time and arrives at the statistical equilibrium state. In
276 such a scenario, the ensemble mean prediction provides no information beyond the mean of
277 the statistical equilibrium state. Therefore, the aim of the long-range forecast is to predict
278 the statistical behavior. In particular, reproducing the statistical equilibrium non-Gaussian
279 probability density function (PDF) with fat tails is a good evidence of successfully predicting
280 the extreme events, where the extreme events and intermittency are the main contributors
281 to the fat tails.

282 Note that different models may have the same characteristics for the long-term statistics
283 but they often have significantly different skill for short and medium range prediction as
284 well as the forced response. In a recent paper⁷⁴, several instructive examples using both a
285 simple linear 2×2 system and more complicated nonlinear models unambiguously illustrate
286 such a feature of predicting complex turbulent dynamical systems. It is also shown in the
287 paper⁷⁴ that in the presence of model error, developing suitable approximate models that
288 are skillful in one of the short-, medium- or long-range forecasting is already a quite difficult
289 task. In many cases, there exists an information barrier⁷⁵ that prevents the approximate
290 models predicting the exact statistics and capturing the perfect response.

291 It is worthwhile to mention that a grand challenge in contemporary climate, atmosphere,
292 and ocean science is to understand and predict intraseasonal variability for time scales from
293 30 to 60 days, which is longer than standard weather time scales of at most a week and
294 much shorter than the yearly time scales of short-term climate. Therefore, it belongs to the
295 medium-range forecasts. Von Neumann⁷⁶ called such problems at the intersection of weather
296 and climate the greatest challenge in future meteorology^{10,77}. The Indian-Asian monsoon
297 and the MJO^{10,78-80} are the most significant intraseasonal variability occurs in the tropical
298 areas. Notably, it is shown in the recent work⁵⁸⁻⁶⁰ that the nonlinear modeling framework
299 in Section II facilitates the development of effective low-order nonlinear stochastic models
300 for predicting the intermittent time series of the MJO and the monsoon as well as extending

301 the predictability of these intraseasonal variabilities.

302 **C. Prediction of the dynamical evolution towards to the attractor**

303 The short-, medium-, and long-range forecasts discussed above typically assume the ini-
304 tial values lie in the statistical equilibrium states. On the other hand, certain internal or
305 external perturbations are able to kick the model variables outside the attractor. Therefore,
306 predicting the time evolution of the extreme events that start from a state outside the at-
307 tractor and its returning path to the statistical equilibrium state is another important issue.
308 Since most approximate models are calibrated using the training data from the attractor,
309 there is no guarantee that these approximate models are automatically able to predict the
310 relaxation towards the attractor. This results in a great challenge of predicting the extreme
311 events starting outside the attractor. Some instructive studies of the prediction and linear
312 response skill with the initial condition being off the attractor can be found in a recent
313 paper⁷⁴.

314 **D. Calibration of the model through parameter estimation**

315 One important issue before applying approximate models for predicting extreme events is
316 the model calibration through parameter estimation. The method adopted here follows the
317 algorithm in a recent work⁸¹. The main difficulty in estimating the parameters in general
318 nonlinear systems with only partial observations is that the closed form of the likelihood
319 function is typically unavailable. Therefore, data augmentation of trajectories associated
320 with the hidden variables is often applied⁸²⁻⁸⁵, which then allows using the Markov Chain
321 Monte Carlo (MCMC) methods to sample the parameters and the hidden trajectories in
322 an alternative way for parameter estimation. Yet, since the hidden trajectories lie in an
323 infinitely dimensional space (or finite but large dimensional with the discrete approximation),
324 the data augmentation can be quite slow in many applications.

Here, in light of the closed analytic formulae (3) of the conditional Gaussian nonlin-
ear approximate models (1), data assimilation can be incorporated into a classical MCMC
algorithm to circumvent the most expensive part of the parameter estimation algorithm,
namely sampling the unobserved trajectories using data augmentation. Specifically, in each

iteration step k of the MCMC, we make use of the observed trajectories \mathbf{u}_I and the current updated parameters $\boldsymbol{\theta}^{(k)}$ to recover the unobserved trajectories of \mathbf{u}_{II} , namely $\mathbf{u}_{II}^{mis,(k)}$ via data assimilation (3) in a deterministic and optimal way. Then $\mathbf{u}_{II}^{mis,(k)}$, \mathbf{u}_I and $\boldsymbol{\theta}^{(k)}$ are used together to compute the likelihood function

$$p(\mathbf{u}_I|\boldsymbol{\theta}^{(k)}) = p(\mathbf{u}_I|\boldsymbol{\theta}^{(k)}; \mathbf{u}_{II}^{mis,(k)}),$$

325 which will be used in the MCMC algorithm for updating the parameters in the $k+1$ iteration
 326 step. For a complete description of the algorithm, see⁸¹ for details. The slight difference of
 327 the algorithm applied here compared to the original version in⁸¹ is that an adaptive MCMC
 328 procedure⁸⁶ for choosing the proposal function is applied.

329 E. Quantifying the prediction skill

330 1. Path-wise measurements

331 The root-mean-square error (RMSE) and the pattern correlation (Corr).

The root-mean-square error (RMSE) and the pattern correlation (Corr) are the two path-wise measurements that have been widely applied to quantify the prediction skill^{21,87-91}. Denote u_i the true signal and \hat{u}_i the prediction estimate, where $i = 1, \dots, n$ is an index in time. These measurements are given by

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (\hat{u}_i - u_i)^2}{n}}, \quad (4)$$

$$\text{Corr} = \frac{\sum_{i=1}^n (\hat{u}_i - \bar{\hat{u}}_i)(u_i - \bar{u}_i)}{\sqrt{\sum_{i=1}^n (\hat{u}_i - \bar{\hat{u}}_i)^2} \sqrt{\sum_{i=1}^n (u_i - \bar{u}_i)^2}}, \quad (5)$$

332 where $\bar{\hat{u}}_i$ and \bar{u}_i denote the mean of \hat{u}_i and u_i respectively.

333 In practice, the trajectory of the ensemble mean is often used as \hat{u}_i in measuring the
 334 RMSE and Corr. These two path-wise measurements are intuitive and easy to be applied.
 335 Typically, a prediction is said to be skillful if the RMSE is below one standard deviation of
 336 the true signal and the Corr is above the threshold $\text{Corr} = 0.5$.

337 Yet, we have to point out several potential issues in these measurements. First, since only
 338 the ensemble mean is used as the predictor, the predicted uncertainty which involves the
 339 ensemble spread (or the confidence interval of the ensemble mean prediction) is not involved

340 in these path-wise measurements. Second, these path-wise measurements fail to quantify
 341 the skill of the long-term forecast, since the ensemble mean simply becomes the equilibrium
 342 mean state of the system. In addition, both the RMSE and Corr take into account only the
 343 information up to the second order statistics. Thus, they may lead to biased conclusions
 344 for predicting extreme events and non-Gaussian features. Nevertheless, due to the simple
 345 form, these path-wise measurements can still be applied to provide some useful information
 346 for the short- and medium-range forecasts.

347 **The temporal autocorrelation function (ACF).**

348 Autocorrelation is the correlation of a signal with a delayed copy of itself, as a function
 349 of delay. For a zero mean and stationary random process u , the autocorrelation function
 350 (ACF) can be calculated as

$$\text{ACF}(t) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \frac{u(t + \tau)u^*(\tau)}{\text{Var}(u)} d\tau, \quad (6)$$

351 where \cdot^* denotes the complex conjugate. The ACF has been widely used to measure the
 352 system memory. It also plays an important role in improving the linear response via the
 353 fluctuation-dissipation theorem^{31,92}. If the perfect model and the approximate model share
 354 the similar ACFs, then the two systems usually have a similar dynamical behavior at least
 355 up to the second order statistics. However, for nonlinear and chaotic systems, high order
 356 statistics may play an important roles for extreme events. Therefore, the ACF can only be
 357 regarded as a crude indicator of the overall predictability of the underlying system. As a
 358 remark, the information theory is able to provide a rigorous and practical way to quantify
 359 the error in the two ACFs associated with the perfect and approximate models by making
 360 use of their spectral representations. See^{30,93} for details.

361 **2. Information measurements**

362 Information theory provides a natural way to quantify the prediction skill and model
 363 error by measuring the lack of information. Different from the path-wise measurements,
 364 the information measurements assess the statistical behavior of the systems. The lack of
 365 information in one probability density q compared with another p is through the relative
 366 entropy $\mathcal{P}(p, q)$ ^{2,94-98},

$$\mathcal{P}(p, q) = \int p \log \left(\frac{p}{q} \right), \quad (7)$$

367 which is also known as the Kullback-Leibler divergence or information divergence^{99–101}. De-
 368 spite the lack of symmetry, the relative entropy has two attractive features. First, $\mathcal{P}(p, q) \geq 0$
 369 with equality if and only if $p = q$. Second, $\mathcal{P}(p, q)$ is invariant under general nonlinear
 370 changes of variables.

371 **Long-term prediction.**

372 The long range forecast using the approximate model aims at capturing the non-Gaussian
 373 statistical equilibrium states of the truth, especially the fat tails that correspond to the ex-
 374 treme events. This is very different from the short- or medium-range forecasts, where the
 375 path-wise measurements of the ensemble mean are informative. In fact, the path-wise mea-
 376 surements completely fail to quantify the long range forecasting skill. Information theory,
 377 on the other hand, provides a natural quantification of the statistical prediction skill in the
 378 approximate model, which is given by

$$\mathcal{E}_{eq} = \mathcal{P}(p_{eq}, p_{eq}^M), \quad (8)$$

379 where p_{eq} and p_{eq}^M are the equilibrium PDFs of the perfect model and the approximate
 380 model, respectively. The information measurement in (8) is able to quantify the skill of the
 381 approximate model in capturing both the majority of the events represented by the mode of
 382 the PDF and the intermittent extreme events in the PDF tails. Note that minimizing the
 383 information score in (8) is also known as capturing the model fidelity⁹⁴ using approximate
 384 models.

385 **The short- and medium-range forecasts.**

386 The information theory can also be applied to quantify the short- and medium-range
 387 forecasting skill. The fundamental difference between the information measurements and
 388 the path-wise ones is that the information measurements are able to take into account
 389 the predicted uncertainty. Denote p_t and p_t^M the PDFs of the time-dependent perfect and
 390 approximate models starting from the same initial time. Similar to (8), an information
 391 metric for quantifying the predicted model error as a function of time can be defined as

$$\mathcal{E}_t = \mathcal{P}(p_t, p_t^M). \quad (9)$$

392 A suitable approximate model is expected to have a small model error throughout the time.

393 The information measurements can also be used to assess the predictability, also known as
 394 the internal prediction skill, of both the perfect and approximate models using the following
 395 matrix^{74,102,103},

$$\mathcal{D}_t = \mathcal{P}(p_t, p_{eq}), \quad \text{and} \quad \mathcal{D}_t^M = \mathcal{P}(p_t^M, p_{eq}^M). \quad (10)$$

396 Clearly, the measurement in (10) quantifies the information provided by the initial condi-
 397 tions about the future state of the system beyond the prior knowledge available through
 398 equilibrium statistics. Obviously, both \mathcal{D}_t and \mathcal{D}_t^M will decay to zero eventually. Therefore,
 399 the measurement in (10) can be regarded as an analog to the ACF but it takes into account
 400 the entire predicted PDF rather than simply the path associated with the ensemble mean
 401 prediction.

402 **IV. A SIMPLE MODEL WITH HIGHLY NON-GAUSSIAN BEHAVIOR IN** 403 **THE HIDDEN PROCESS**

404 Stochastic parameterizations are widely used in developing approximate models for com-
 405 plex dynamical systems with partial observations^{1,104–106}. The idea of applying stochastic
 406 parameterizations is to use simple stochastic processes to describe the complicated dynam-
 407 ics of the unobserved or unresolved scales such that the overall computational cost of the
 408 approximate models is greatly reduced. One important practical issue is to develop suitable
 409 stochastic parameterizations for the hidden processes such that the intermittent features
 410 are captured and the approximate models with the stochastic parameterizations are able to
 411 accurately predict the extreme events in the observed variables.

412 The goal of this section is to test the skill of a simple and efficient stochastic parameter-
 413 ization strategy in predicting intermittent non-Gaussian features and extreme events based
 414 on a low-order highly non-Gaussian test model given only a short period of training data
 415 with partial observations.

416 **A. The perfect and approximate models**

417 **The perfect model.**

The perfect test model here is given by a two-dimensional system where only a short

trajectory of one variable u is observed. The model reads,

$$du = \left(-\gamma u + F_u \right) dt + \sigma_u dW_u, \quad (11a)$$

$$d\gamma = (a_\gamma \gamma + b_\gamma \gamma^2 + c_\gamma \gamma^3 + f_\gamma) dt + (A_\gamma + B_\gamma \gamma) dW_{\gamma,1} + \sigma_\gamma dW_{\gamma,2}. \quad (11b)$$

418 In this model, the variable γ acts as a stochastic damping in the equation of u and the aver-
 419 aged value of γ over time needs to be positive to guarantee the mean stability of u ¹⁰⁷. Once
 420 the sign of γ switches from positive values to negative values, γ becomes anti-damping and
 421 it leads to the intermittent events in u . On the other hand, γ is driven by a cubic nonlinear
 422 equation with correlated additive and multiplicative noise. This cubic model is a canon-
 423 ical model for low frequency atmospheric variability^{108,109}. This one-dimensional, normal
 424 form has been applied in a regression strategy for data from a prototype atmosphere and
 425 ocean model to build one-dimensional stochastic models for low-frequency patterns such as
 426 the North Atlantic Oscillation and the leading principal component that has features of the
 427 Arctic Oscillation. Given the non-Gaussian features and the potential physical explanations,
 428 the low-order model (11) becomes a useful testbed for developing suitable stochastic param-
 429 eterization strategies of the hidden process that allows skillful prediction of the extreme
 430 events in the observed variable.

431 The following parameters are taken in the perfect model,

$$\begin{aligned} F_u = 0.3, \quad \sigma_u = 0.1, \quad a_\gamma = -\frac{3}{8}, \quad b_\gamma = 1, \quad c_\gamma = -\frac{1}{2}, \\ A_\gamma = 0, \quad B_\gamma = \frac{1}{2\sqrt{2}}, \quad f_\gamma = 0.1, \quad \sigma_\gamma = \frac{1}{2\sqrt{2}}. \end{aligned} \quad (12)$$

432 With these parameters, the model trajectories together with the equilibrium PDFs and ACFs
 433 are shown in Panels (a)–(c) of Figure 1. Note that the time series in Panel (a) only contains
 434 a length of 500 time units but the PDFs and ACFs in Panels (b)–(c) are computed based
 435 on the model simulation with a length of 10,000 units in order to minimize the sampling
 436 bias in showing these statistics.

437 In this dynamical regime, the time series of γ shows a stochastic switching behavior.
 438 Roughly speaking, γ has two statistical states. The averaged value in one state is slightly
 439 negative, corresponding to the intermittent phase of u , while another state of γ is positive,
 440 corresponding to the quiescent phase of u . The PDF of u , due to the intermittent extreme
 441 events, is highly skewed with an one-sided fat tail. On the other hand, the PDF of γ shows a

442 bimodal behavior, which is also significantly non-Gaussian. The ACFs indicate that overall
 443 u has a longer memory than γ .

444 **The approximate model.**

445 The perfect model (11) here can be regarded as a paradigm model in many real applica-
 446 tions, where the hidden variables are driven by some unknown complicated processes that
 447 interact with the observed variables in a highly nonlinear way. From a practical point of
 448 view, it is important to develop a simple and computationally tractable approximate mod-
 449 el which is nevertheless able to capture the key nonlinear feedback from the unobserved
 450 variable γ to the observed variable u . The approximate model is expected to predict the
 451 extreme events of the observed process u .

452 One commonly used reduced order modeling strategy is to adopt a mean stochastic model
 453 (MSM) for the observed process u . The MSM makes use of the averaged value of γ as the
 454 damping term and the resulting system is

$$du = (-\hat{\gamma}u + F_u)dt + \sigma_u dW_u. \quad (13)$$

455 Since the mean stability is guaranteed in the original system, the constant $\hat{\gamma}$ is positive.
 456 Thus, the MSM is a linear model with Gaussian statistics. It has been shown in^{22,107} that
 457 the MSM is unable to capture the short-term rapid increment of the intermittent trajectory
 458 of u due to the lack of intermittent instability mechanism. Such a Gaussian model also fails
 459 to predict the long-term non-Gaussian PDF with skewness and fat tails.

Here, a new approximate model is developed using the stochastic parameterized equation
 technique^{48,49}, the idea of which has been applied to the extended Kalman filters (known as
 the SPEKF-type model) and other prediction and data assimilation forecast models. The
 approximate model has the following form,

$$du = (-\gamma u + F_u)dt + \sigma_u dW_u, \quad (14a)$$

$$d\gamma = -d_\gamma(\gamma - \hat{\gamma})dt + \sigma_\gamma dW_\gamma. \quad (14b)$$

460 In (14), the nonlinear process γ with correlated additive and multiplicative noise in (11b)
 461 has been simplified to a linear process with only Gaussian additive noise. Nevertheless, the
 462 variable γ remains switching between positive and negative phases, representing damping
 463 and anti-damping effects as a feedback to u . Therefore, the variable γ is still able to trigger

464 intermittent extreme events in u . One important feature is that the approximate model
 465 (14) belongs to the conditional Gaussian nonlinear framework as was described in Section
 466 II, which allows the effective algorithm (3) to solve the conditional statistics of the hidden
 467 variable γ given the observations from u . This greatly facilitates the data assimilation and
 468 predictions.

469 B. Parameter estimation

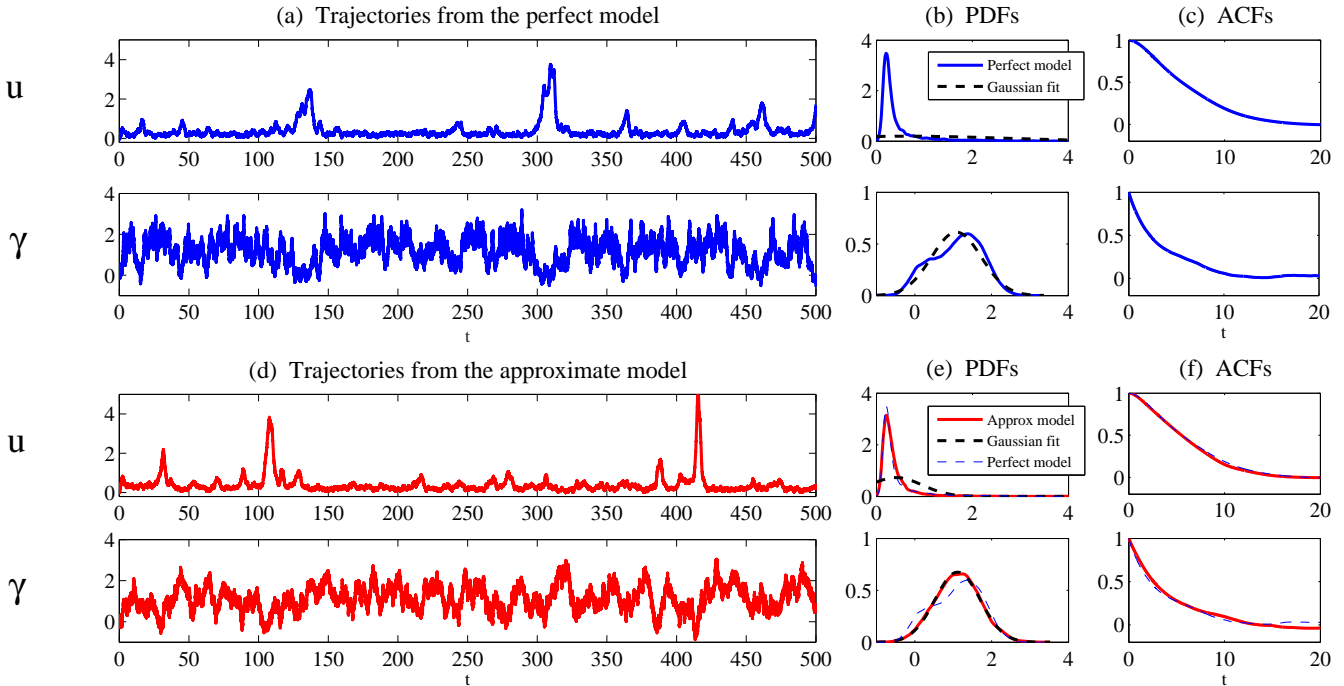
470 Before applying the approximate model for prediction, the parameters $(F_u, \sigma_u, d_\gamma, \hat{\gamma}, \sigma_\gamma)$
 471 in the approximate model (14) need to be estimated. The training time series only involves
 472 the observed variable u and the training data has only a short period with 500 units as
 473 shown in Panel (a) of Figure 1. Applying the parameter estimation algorithm described
 474 in Section III D, the results are shown in Figure 2. The trace plots associated with all the
 475 parameters clearly indicate the convergence towards certain values with small uncertainties.
 476 Notably, the estimation values of the two parameters σ_u and F_u in the observed process are
 477 almost the same as the ones in the perfect model. The averaged values of the trace plots
 478 from iteration $k = 5000$ to iteration $k = 10000$ are utilized as the estimated parameters in
 479 the approximate model for prediction:

$$d_\gamma = 0.2545, \quad \hat{\gamma} = 1.121, \quad F_u = 0.2489, \quad \sigma_\gamma = 0.4362, \quad \sigma_u = 0.1008. \quad (15)$$

480 C. Long-term prediction

481 With the estimated parameters, we begin with studying the long-term prediction using
 482 the approximate model (14). As an analogy to Panels (a)–(c) in Figure 1 for the perfect
 483 model, Panels (d)–(f) in Figure 1 show the trajectories, PDFs and ACFs of the approximate
 484 model. Note that Panel (d) is simply a free run of the model. Therefore, there is no point-
 485 to-point correspondence between the trajectories shown in Panels (a) and (d) for the perfect
 486 and approximate models. Nevertheless, it is easy to see that the trajectories from the perfect
 487 and approximate models are qualitatively similar to each other, indicating the skill of the
 488 approximate model in capturing the long-term dynamical and statistical behavior. Next, to
 489 understand the quantitative similarity between the two models, the equilibrium PDFs and
 490 the ACFs are compared.

491 Panels (c) and (f) show that both the ACFs of u and γ associated with the two models
 492 are very similar to each other, indicating the success of the approximate model in capturing
 493 the temporal information of the perfect system. On the other hand, as shown in Panels
 494 (b) and (e), the PDF of u is also perfectly recovered by the approximate model, where
 495 using the information distance (7) the difference between the PDFs associated with the
 496 approximate and perfect models $\mathcal{P}(p_{eq}(u), p_{eq}^M(u)) = 0.0345$ is a negligible value. The PDF
 497 of γ is not perfectly recovered because the approximate model uses only a linear system with
 498 additive noise for γ , which fails to capture the non-Gaussian PDF of γ . This is known as
 499 the information barrier²². Nevertheless, the PDF of γ associated the approximate model is
 500 nearly exact the same as the Gaussian fit of the bimodal distribution associated with the
 501 truth. This in fact implies that the approximate model has reached its predictability limit
 502 in predicting the long-term statistics of the hidden variable γ .



503

504 FIG. 1. Model trajectories, PDFs and ACFs of the perfect model (11) with parameters in (12).
 505 Top: u ; Bottom: γ . Note that the PDFs and ACFs are computed based on the model simulation
 506 with a length of 10,000 units. But in panels (a) and (c) only time series with a length of 500 time
 507 units are shown. The black curves in the PDFs show the Gaussian fits.

508

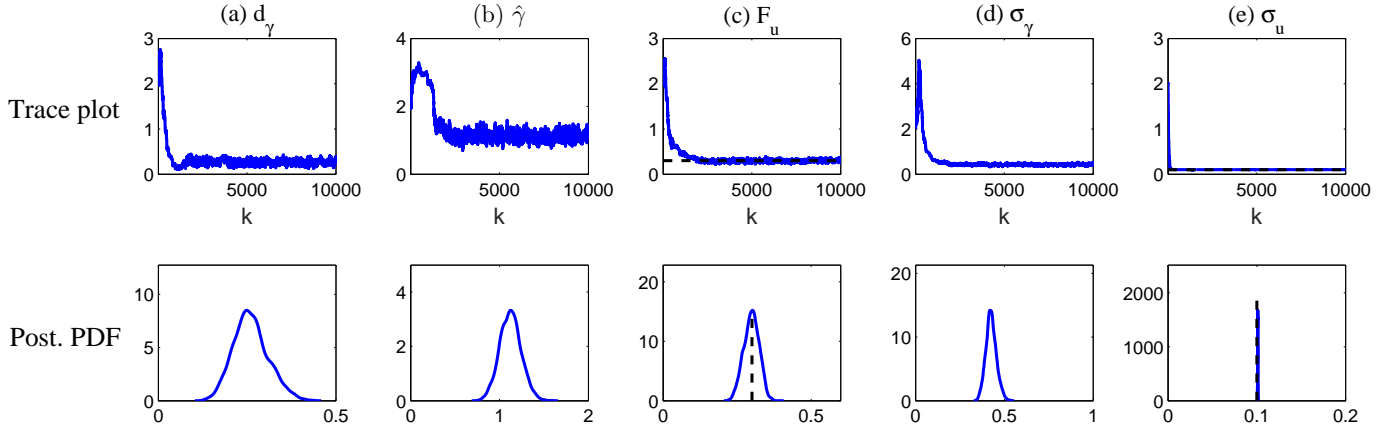


FIG. 2. Parameter estimation of the approximate model (14). Top: trace plot. Bottom: posterior PDFs of the parameters from the trace plot taking the values from $k = 5000$ to $k = 10000$. The black lines show the values of σ_u and F_u in the perfect model (11), serving as reference values.

510 D. Data assimilation

511 One key feature of the approximate model (14) is that it belongs to the conditional
512 Gaussian model family (1), which allows using the closed analytic formulae (3) to solve the
513 the conditional distribution $p(\gamma|u)$ for assimilating the unobserved variable γ . Note that
514 the perfect model (11) is not a conditional Gaussian system and expensive particle methods
515 have to be used in order to assimilate the unobserved variable γ even in this two-dimensional
516 system. Therefore, the approximate model (14) is much more computationally efficient for
517 state estimation, data assimilation and prediction.

518 Figure 3 shows the data assimilation results using the approximate model (14) as the
519 forecast model. It is clear that the γ values associated with the intermittent phase of u
520 are recovered with both high accuracy and low uncertainty. The accurate recovery of the
521 hidden variable γ at the intermittent phase of u indicates its potential for predicting the
522 extreme events. On the other hand, assimilating the γ states corresponding to the quiescent
523 phase of u are recovered with high uncertainty. The posterior mean also fails to track the
524 fluctuations in the true trajectory. This is not surprising since the quiescent phases of u
525 have weak amplitudes and therefore the noise-to-signal ratio is large. In fact, as long as the
526 hidden variable γ stays positive, playing the role as a damping, it has very weak influence
527 on the dynamics u at the quiescent phases. The assimilated values and uncertainties of γ
528 accurately reflect these features.

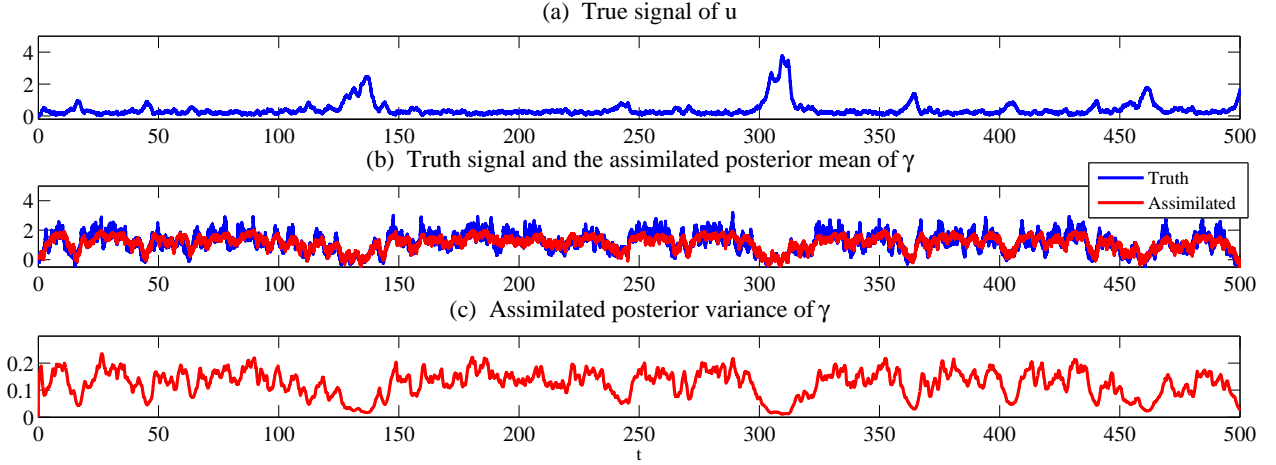


FIG. 3. Data assimilation of the hidden variable γ using the approximate model (14) as the forecast model. The true signal of the observed variable (panel (a)) is generated from (11). Panel (b) shows the true signal of γ from (11) and the assimilated (filtered) posterior mean of γ using the approximate forecast model (14). Panel (c) shows the posterior variance. The black dashed boxes indicates the events that will be studied for short-term prediction in the next subsections.

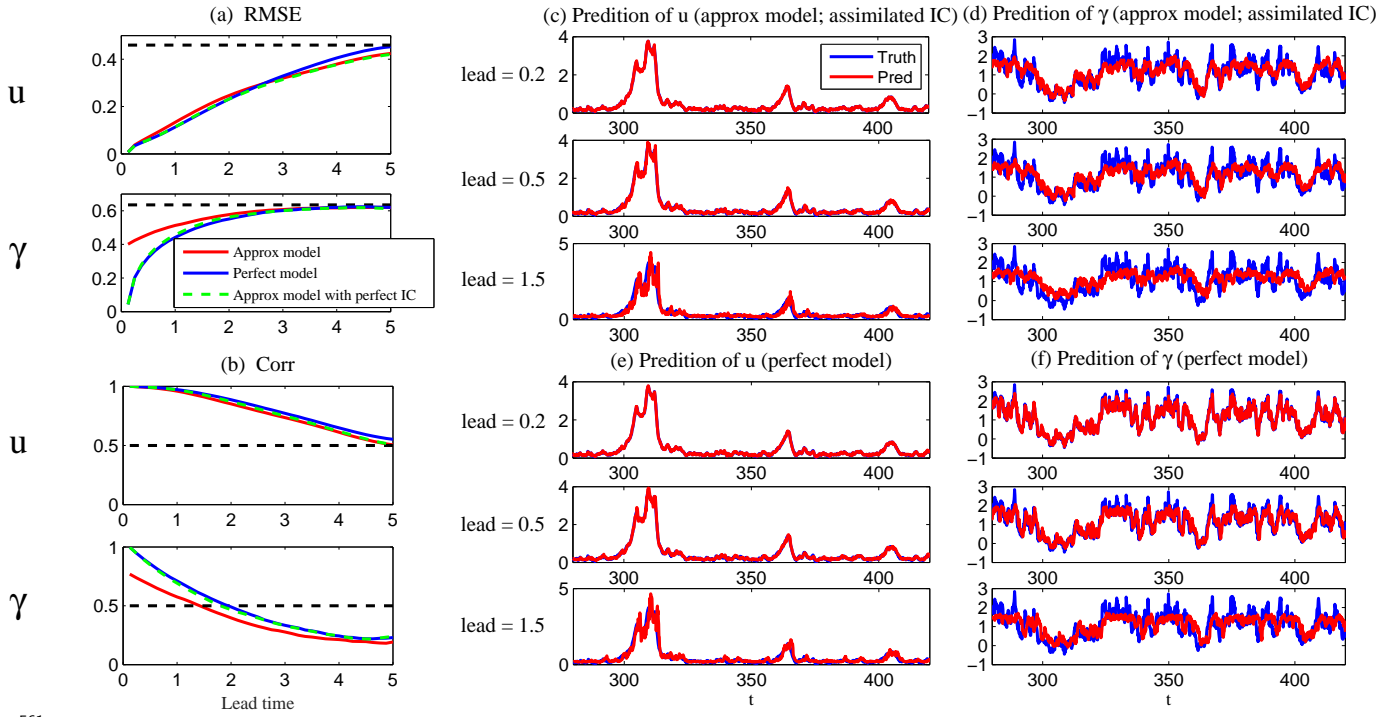
529 E. Short- and medium- range forecasts

530 To study the short- and medium-range forecast, we first show the RMSE and the Corr
531 between the predicted time series and the truth as a function of lead time. Here the ensemble
532 mean is used as the predicted time series. As illustrated in Panels (a)–(b) of Figure 4,
533 the approximate model with the assimilated initial conditions has an overall comparable
534 prediction skill as the perfect model prediction with the perfect initial conditions. The only
535 main difference lies in the very short term for predicting γ , where the prediction using the
536 approximate model with the assimilated initial conditions has a larger error. This is due
537 to the large uncertainty in the assimilated initial conditions at the quiescent phases. In
538 fact, if we adopt the approximate model as the forecast model but use the perfect initial
539 conditions (green curves), then the prediction skill is almost the same as using the perfect
540 model prediction. Note that the overall skillful prediction of u lasts up to 5 units while that
541 of γ is around 2 units.

542 Panels (c)–(d) and (e)–(f) of Figure 4 show the lead time prediction at 0.2, 0.5 and 1.5
543 units using the approximate model with the assimilated initial conditions and the perfect
544 model with the perfect initial conditions, respectively. The prediction of u , especially the
545 extreme events, is quite accurate at all the three lead times for both the models. The
546 prediction of the negative phase of γ is also nearly perfect. The only difference between the
547 two models lies in predicting the positive phases of γ , where the approximate model cannot
548 provide an accurate prediction even at a very short lead time. This is due to the error and
549 the uncertainty in the assimilated initial conditions as was discussed above. On the other
550 hand, while the perfect model is able to predict the positive phase of γ (corresponding to the
551 quiescent phases of u) in a very short term, it is interesting to see that even with the perfect
552 model and perfect initial conditions, some significant errors already appear in predicting the
553 positive phases of γ at a lead time 0.5. At a lead time 1.5, the perfect model essentially gives
554 the same results as the approximate model, where an accurate prediction is found in both
555 u and the negative phase of γ while the model is not very skillful in predicting the positive
556 phase of γ . These facts indicate that when γ is positive it only has a weak influence on u
557 and therefore the system has an intrinsic weak dependence of γ .

558 To conclude, the approximate model has almost the same short- and medium-range fore-
559 casting skill as the perfect model, especially in predicting the extreme events in u and the

560 corresponding triggering phases in γ .



561

562 FIG. 4. Short- and medium-range forecasts. Panels (a)–(b): RMS error and pattern correlation
 563 between the predicted time series and the truth as a function of lead time. Red: prediction
 564 using the approximate model (14), where the initial values of γ are obtained by data assimilation.
 565 Dashed blue: prediction using the perfect model (11) with perfect initial conditions. Dashed green:
 566 prediction using the approximate model (14) but with perfect initial conditions. Panels (c)–(d):
 567 Ensemble mean prediction using the approximate model with assimilated initial condition (IC) at
 568 lead times 0.2, 0.5 and 1.5. The blue curves show the truth while the red ones show the prediction.
 569 Panels (e)–(f): Similar to (c)–(d) but using the perfect model and perfect initial condition.
 570

571

572 **F. Prediction with an initial value starting outside the attractor**

573 Finally, we study the prediction skill of the approximate model if the initial value is out-
574 side the attractor (the statistical equilibrium state). In Figure 5, we consider the situations
575 where either the initial value of u or that of γ is outside the attractor. It is clear that when
576 $u(0)$ is outside the attractor while γ stays in the attractor (Panels (a) and (c)), the trajectory
577 of u releases to the attractor in a similar fashion using both the approximate model and
578 the perfect model. This is because there is no approximation in the observed process u and
579 the time evolution of γ at the attractor has already been shown to be accurately described
580 using the approximate model. On the other hand, if γ starts from a value that is outside the
581 attractor (Panels (b) and (d)), then the approximate model in capturing the relaxation of γ
582 towards the attractor may contain errors. In fact, when γ starts from a large value as shown
583 in Panel (b), the cubic damping plays an important role in strongly pushing the system
584 towards the attractor. Starting from a large value of γ , the impact of the cubic damping is
585 much stronger than that at the attractor and therefore the approximate model with a linear
586 damping in γ fails to capture this feature. Nevertheless, if u stays on the attractor, then as
587 long as γ is positive it has only a weak influence on the observed variable u . Therefore, the
589 overall dynamics of u can still be described quite well using the approximate model.

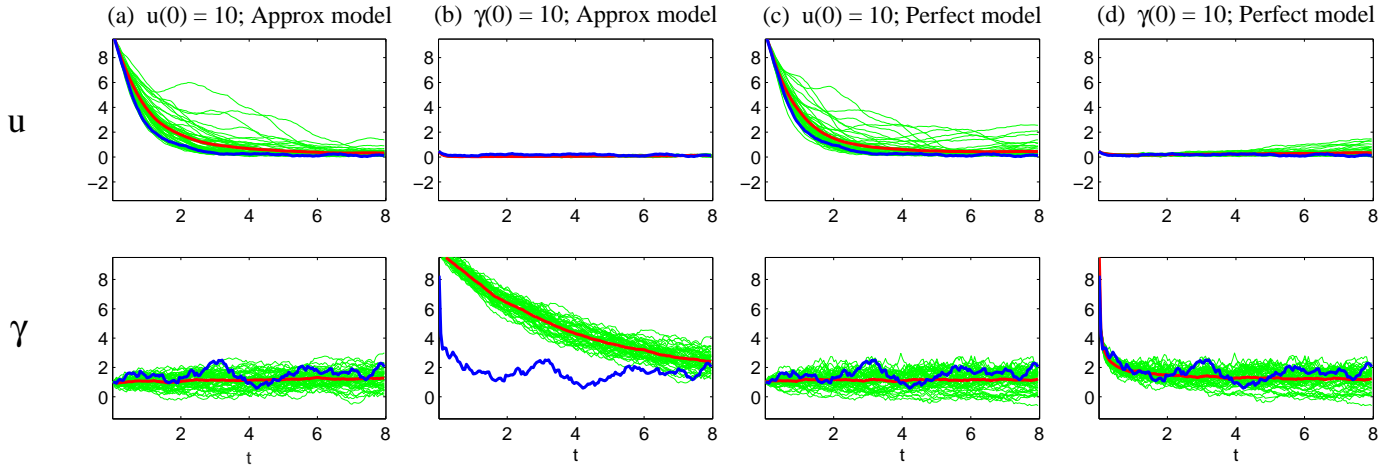


FIG. 5. Ensemble predictions with the initial values starting outside the attractor. Panels (a) and (c): $u(0) = 10$ starts from a value that is off the attractor. $\gamma(0) = 1$ is inside the attractor. Panels (b) and (d): $u(0) = 0.5$ starts from a value that is inside the attractor. $\gamma(0) = 10$ is off the attractor. Here Panels (a)–(b) show the results using the approximate model while Panels (c)–(d) show those using the perfect model. In all the panels, blue curves show the truth and red curves show the ensemble mean which is the average value of 50 ensembles showing in green color.

590 **V. A DYAD MODEL WITH ENERGY-CONSERVING NONLINEAR**
 591 **INTERACTION**

592 The nonlinear test model in the previous section involves only an one-way influence from
 593 γ to u . Yet, in many applications, the observed variables and the unobserved ones have
 594 mutual interactions, which are also often through energy-conserving nonlinear terms^{35,36}.
 595 Therefore, it is important to understand different strategies in building approximate models
 596 to predict the extreme events and other non-Gaussian behavior in such kind of the systems.
 597 In this section, a simple but judicious feedback control strategy is adopted to facilitate the
 598 prediction of the hidden extreme events in an energy-conserving nonlinear dyad model.

599 **A. The models**

600 **The perfect model.**

601 Consider a nonlinear dyad model with energy-conserving nonlinear interaction,

$$\begin{aligned} dv &= \left(-d_v v - cu^2 \right) dt + \sigma_v dW_v, \\ du &= \left((-d_u + cv)u + F_u \right) dt + \sigma_u dW_u. \end{aligned} \tag{16}$$

602 Again only partial observations are available in this nonlinear dyad model, where v is the
 603 observed variable while u is unobserved. This low-order nonlinear model can be regarded
 604 as a toy model of complex turbulent flows. For example, v can be treated as one of the
 605 Fourier modes associated with the large-scale observed variables while u is associated with
 606 the hidden mechanism that drives v . If u represents unresolved or small-scale variables, then
 607 its statistics can be highly non-Gaussian. Here, v plays the role of the stochastic damping
 608 in the process of u such that intermittent extreme events appear in the trajectory of u .
 609 Note that this model is very different from the SPEKF-type of the model described in the
 610 previous section. In fact, in the dyad model (16), the variable u also provides a nonlinear
 611 feedback to v via $-cu^2$ such that the total energy in the nonlinear terms of the coupled
 612 system is conserved, which is known as the physics constraint^{35,36}.

613 Below, the nonlinear dyad model (16) is used as the perfect model. The focus of this
 614 section is to predict the extreme events in the unobserved process u . To this end, the

615 following parameters are taken in the nonlinear dyad model (16),

$$F_u = 1, \quad d_u = 0.8, \quad d_v = 0.8, \quad \sigma_u = 0.2, \quad \sigma_v = 2, \quad c = 1.2. \quad (17)$$

616 As shown in Panels (a)–(c) of Figure 6, the nearly Gaussian observed variable v switches
617 between positive and negative states, which leads to the intermittency in the hidden process
618 u . The non-zero forcing $F_u = 1$ makes the signal of u stay almost within the positive values
619 and the PDF of u is skewed with an one-side fat tail. Note that the amplitude of this forcing
620 term provides different dynamical behavior of the model. In the last part of this section, the
621 prediction skill in different dynamical regimes with various values of F_u will be reported.

622 **The approximate model.**

623 Again, a suitable approximate model is able to predict the extreme events and other
624 important non-Gaussian features of the perfect model. Meanwhile, the approximate model
625 is expected to be computationally efficient for data assimilation and prediction. Due to
626 the closed analytic formulae of the conditional Gaussian models in assimilating the unob-
627 served variables, we now aim at developing a suitable approximate model that belongs to
628 the conditional Gaussian framework. Note that by observing v , the perfect model (16) is
629 not a conditional Gaussian nonlinear system. One starting idea for building an approximate
630 model is to apply a bare truncation strategy, which ignores the quadratic feedback term $-u^2$
631 in the process of v in (16). This is actually a commonly used strategy in developing approx-
632 imate models for many complicated systems in practice, where some nonlinear terms are
633 dropped. However, this strategy does not work for studying the extreme events with partial
634 observations. In fact, without this feedback term, the variable u is completely decoupled
635 from the process of v . In other words, given only the observations in v , the processes and
636 the parameters of u are not even identifiable. What is more, using the same parameters as
637 in (17), such an approximate model suffers from a finite-time blowup of the signals^{35,110}.

638 The failure of the bare truncation model is due to the complete ignorance of the nonlinear
639 feedback term from u to v . This nonlinear feedback not only provides the observability of
640 u in the v process but also offers the important causal effects between the two processes.
641 Therefore, a suitable approximate model is supposed to take into account such an interaction

642 between the two processes. To this end, the following approximate model is adopted,

$$\begin{aligned}
 dv &= \left(-d_v v - cu \right) dt + \sigma_v dW_v, \\
 du &= \left((-d_u + cv)u + F_u \right) dt + \sigma_u dW_u.
 \end{aligned}
 \tag{18}$$

643 This approximate model uses a linear feedback $-cu$ to approximate the nonlinear interaction
 644 $-cu^2$ in the original dyad model. This simplification can be regarded as using a linear control
 645 term to retain the mutual dependence of u and v . It also allows the approximate model to
 646 belong to the conditional Gaussian framework that facilitates efficient data assimilation and
 647 prediction algorithms.

648 B. Parameter estimation

649 For the parameter estimation of the approximate model (18), we make use of a short
 650 training data of v with only 500 time units as shown in Panel (a) of Figure 6. The parameter
 651 estimation algorithm is run for $K = 15000$ steps and the averaged values from the trace
 652 plots between $k = 5000$ to $k = 15000$ is used as the estimated parameters,

$$\begin{aligned}
 d_v &= 0.9234, & d_u &= 0.6672, & c &= 1.8249, & F_u &= 0.6041, \\
 \sigma_u &= 0.0527, & \sigma_v &= 2.0203.
 \end{aligned}
 \tag{19}$$

653 It is useful to compare the estimated parameter values in the approximate model (19) with
 654 those in the perfect model (17). This helps understand the dynamical properties of the
 655 approximate model.

656 The feedback parameter c in the approximate model (17) is increased. This is due to
 657 the fact that cu^2 in the perfect model is replaced by cu in the approximate model while
 658 the amplitude of u in the perfect model is often larger than 1 especially in the intermittent
 659 phases. Therefore, the coefficient c has to be increased in order to retain the amplitude of
 660 the feedback from u to v . On the other hand, according to the second equation in (18), due
 661 to the increase of c , the amplitude of u will increase as well especially for the intermittent
 662 phase. Therefore, the forcing F_u in the approximate model is decreased in order to retain
 663 the amplitude of the observed variable u as in the perfect model.

664 **C. Long-term prediction**

665 With the estimated parameters in hand, we first compare the long range forecasts between
 666 the perfect dyad model (16) and the approximate model with the linear feedback (18).

667 In Panels (d)–(f) of Figure 6, the trajectories, the PDFs and the ACFs associated with
 668 the approximate model are shown, where for a fair comparison of the time series, the same
 669 random number seeds are used. The recovered trajectory of the observed variable v using
 670 the approximate model with the linear feedback term almost perfectly matches that of the
 671 truth (with $\text{Corr} = 0.998$ and $\text{RMSE} = 0.011$).

672 Now let us focus on the hidden intermittent variable u . Comparing the second and the
 673 fourth rows of Figure 6, it is clear that the approximate model with the linear feedback
 674 (18) is skillful in generating the intermittent extreme events in u . In fact, the pattern
 675 correlation between the two time series in these two rows is 0.93, which also indicates that
 676 the approximate model is able to capture the timing of the occurrence of extreme events.
 677 Yet, there are two main errors in the approximate model. First, the amplitudes of the
 678 intermittent events seem to be slightly overestimated. This is easy to understand because in
 679 order to reach the same observed trajectory v , the linear feedback requires a larger u in the
 680 approximate model than the quadratic nonlinear feedback in the perfect model. Second, the
 681 quiescent events also seem to be slightly underestimated in the approximate model. This
 682 results in the fact that the peak of the associated PDF is closer to zero than that of the true
 683 signal. The model error

$$\mathcal{E}_{eq} = \mathcal{P}(p_{eq}, p_{eq}^M) = 0.46, \quad (20)$$

684 which, although is non-negligible, comes largely from the quiescent events. The PDFs of
 685 u associated with both the perfect and approximate models are skewed with an one-sided
 686 fat tail. Therefore, the long range behavior of the approximate model in capturing the
 687 information in the tail that corresponds to the extreme events remains similar to that of the
 688 perfect model.

689 Another conclusion drawn from Figure 6 is that the ACFs of u and v associated with
 690 both the perfect model and the approximate model are very similar to each other, decaying
 691 to zero after one time unit.

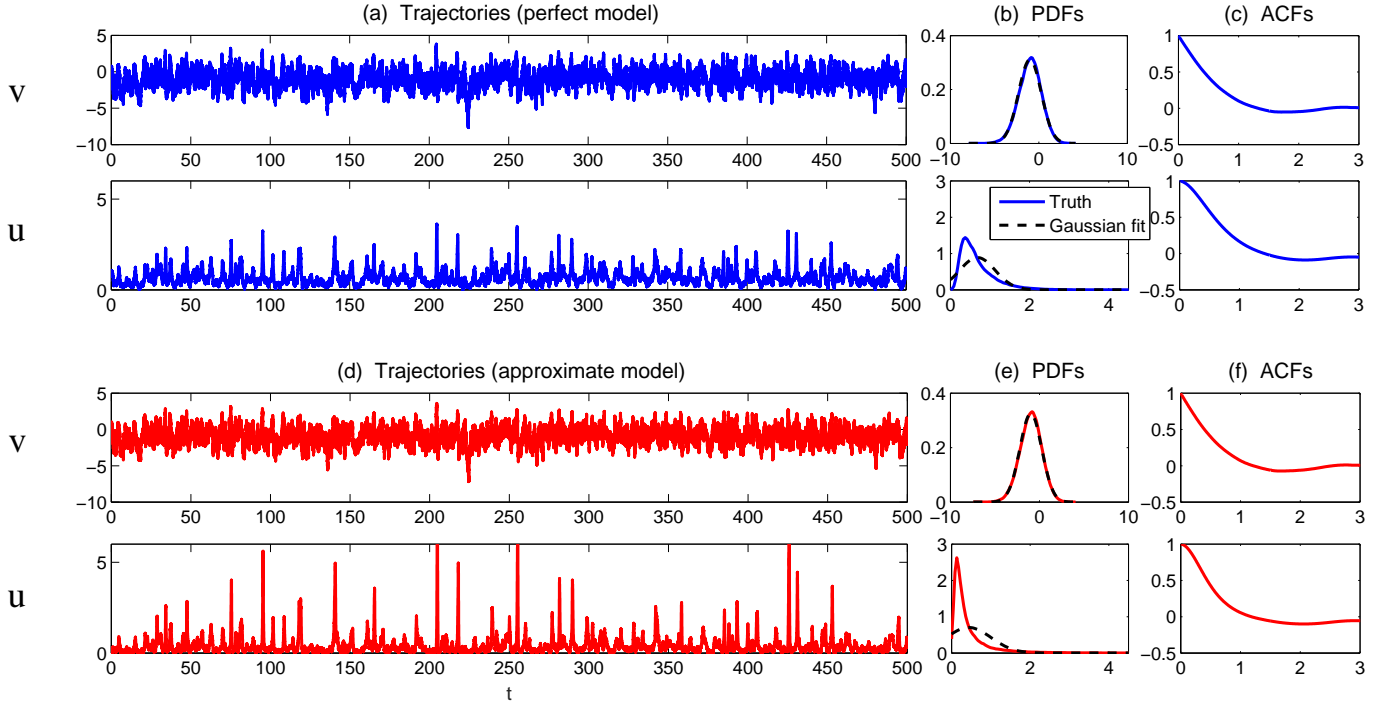


FIG. 6. Panels (a)–(c): Time series, PDFs and ACFs of the dyad model (16) with the parameters in (17). Panels (d)–(f): those of the approximate model (17) with the estimated parameters in (19).

693 D. Data assimilation, Short- and medium-range forecasts

694 Given the observation in v , the assimilated u is shown in Figure 7. Overall, the assimilated
695 signal of the hidden variable u and the truth have a very good match in terms of the patterns.
696 Yet, due to the intrinsic model error as discussed above, the quiescent and intermittent
697 phases are slightly underestimated and overestimated, respectively.

698 Panels (a)–(b) of Figure 8 show the RMSE and the Corr between the true signal and the
699 ensemble mean predictions as a function of lead time. Except at the very short lead time,
700 where the data assimilation results in some uncertainties in the initial values, the approxi-
701 mate model essentially gives the same prediction skill as the perfect model in terms of the
702 RMSE and the Corr. This indicates the overall skillful prediction using the approximate
703 model. Note that since our focus is the extreme events in the hidden process, some extra
704 information beyond the RMSE and Corr needs to be explored. In Panels (c)–(d), a compar-
705 ison of the medium range forecasts and the forecast PDFs at lead time $t = 0.6$ is shown. It
706 is clear that the approximate model is more skillful in capturing the extreme events and the
707 fat tail of the predicted PDF than the perfect model. This is not surprising. In fact, it is
708 well known that the amplitude of the ensemble mean prediction decays as time evolves. On
709 the other hand, the slight overestimation of the amplitude of u in the approximate model
710 compensates the underestimation of the amplitudes in the ensemble mean forecast, which is
711 crucial in predicting extreme events at the medium range.

712 Figure 9 shows four case studies of the time evolution of the predicted PDFs starting
713 from different initial phases. The predictions of v using both the perfect model and the
714 approximate model are overall similar to each other. Note that more ensemble members in
715 the prediction using the approximate model are actually able to forecast the extreme events
716 than the perfect model. This feature is quite useful for medium-range forecast, especially
717 when the starting time is an onset phase of the extreme events in u (Cases 1 and 2). In
718 addition, even with some errors in the initial condition due to the data assimilation (Case
719 4), the approximate model is still able to capture the time evolution of extreme events.

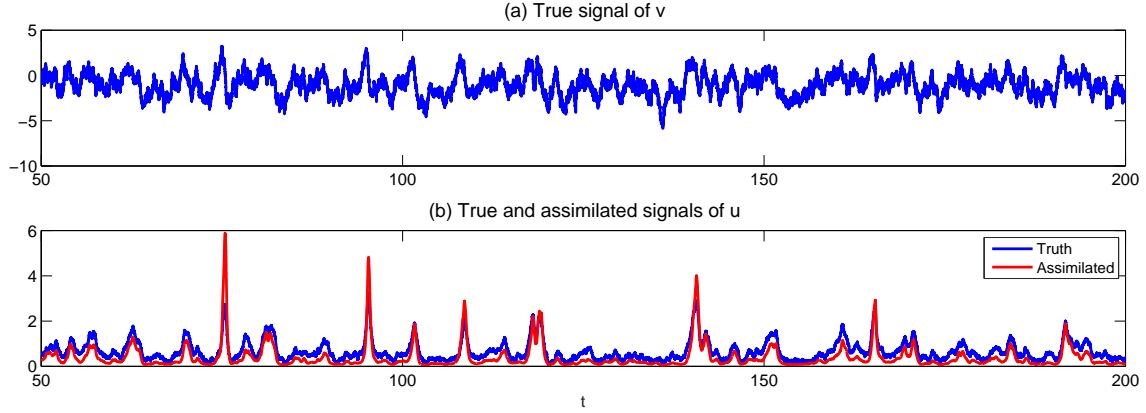


FIG. 7. Panel (a): the true signal of v from the perfect dyad model (16). Panel (b): the true signal of u from the perfect dyad model (16) and assimilated signals of u using the approximate model (18).

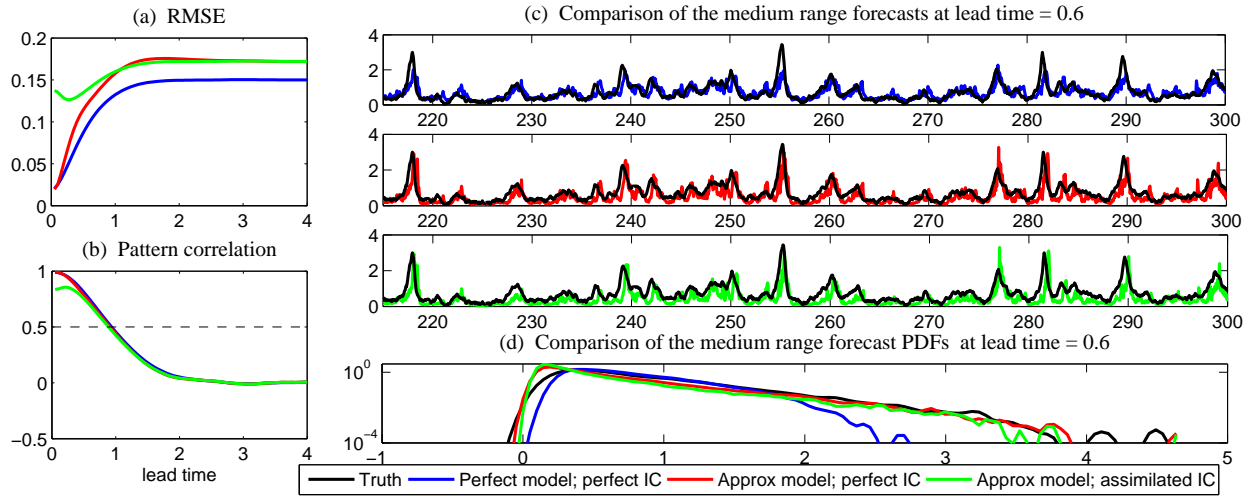


FIG. 8. Short- and medium-range forecasts based on the ensemble mean forecast. Panel (a)–(b): RMSE and Corr between the true signal and the prediction ones as a function of lead time. Blue: perfect model prediction with the perfect initial condition. Red: approximate model prediction with perfect initial condition. Green: approximate model prediction with assimilated initial conditions. Panel (c): comparison of the medium range forecasts at lead time $t = 0.6$. Panel (d): comparison of the medium range forecast PDFs at lead time $t = 0.6$.

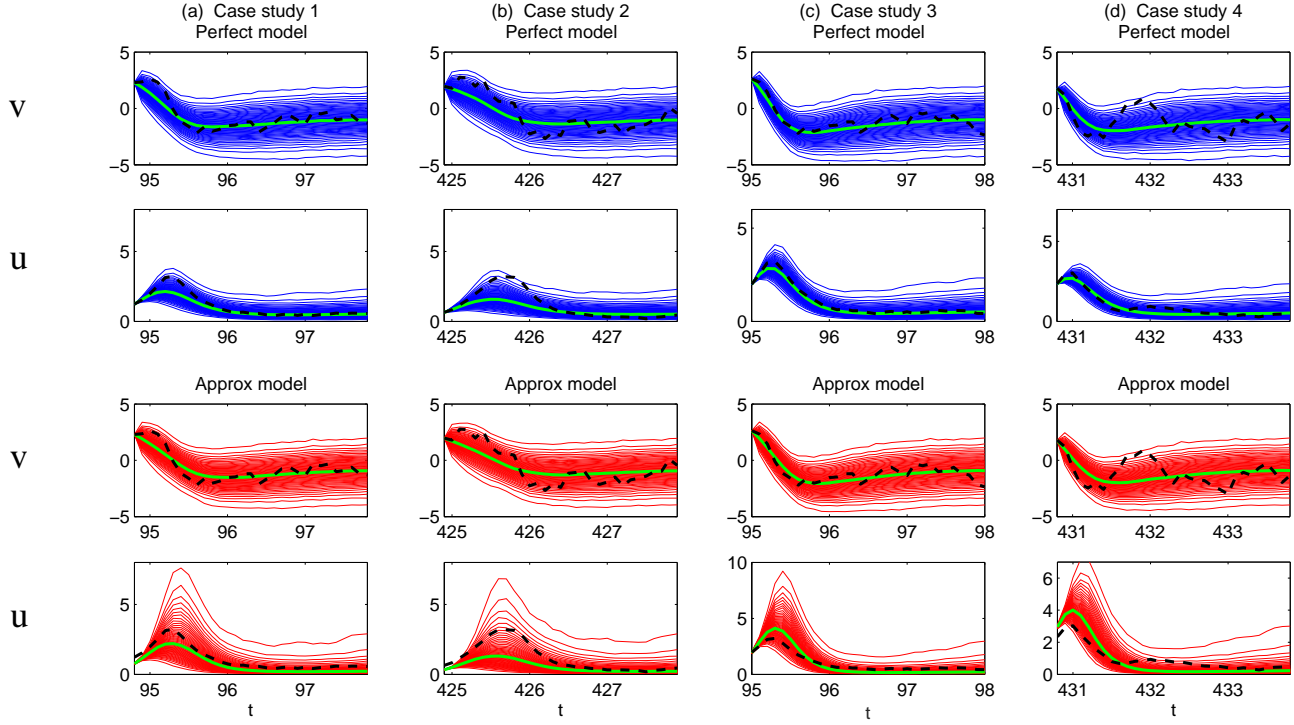


FIG. 9. Case studies. Time evolution of the predicted PDFs starting from different initial phases. Each PDF is shown with 50 thin curves, which represent the 1st, 3rd, 5th, ..., 97th and 99th percentiles of the of the PDF. The green curve represents the mode of the PDF since the PDF is non-Gaussian. The black dashed curve is the true signal. Columns (a)–(b): starting from an onset phase of extreme events. Columns (c)–(d): starting from a mature phase of extreme events.

720 E. Prediction with an initial value starting outside the attractor

721 In this subsection, we compare the predictions when the initial values are outside the
722 attractor. According to Figure 6, the attractor of u contains values that are positive but
723 almost always stay below $u = 4$. Thus, we consider the following two situations: A). the
724 hidden variable u starting from a negative value, and B). the hidden variable u starting from
725 a large positive value.

726 A. The hidden variable u starting from a negative value.

727 In Panels (a)–(f) of Figure 10, we show the prediction where u starts from a negative
728 value $u(0) < 0$. Here, v always starts from its equilibrium mean value $v(0) = -0.9584$.
729 In Panels (a)–(b), $u(0) = -0.2$ is slightly negative. The approximate model behaves in a
730 similar way as the perfect model, where after a short term, the trajectory will arrive at the
731 attractor. However, when $u(0) = -0.5$ as shown in Panels (c)–(d), some of the ensemble
732 members in the approximate model blows up in finite time (around $t = 1.5$). See the second
733 row of Panel (d). Such a behavior becomes even worse when $u(0)$ is decreased to $u(0) = -0.8$
734 as shown in Panels (e)–(f), where quite a few ensemble members blow up in a short finite
735 time (around $t = 0.5$ to $t = 1.5$). Panel (g) of Figure 10 shows the percentage of the events
736 that blow up as a function of the initial value $u(0)$. As expected, with the decrease of $u(0)$,
737 the number of blowup events increases.

738 Now we look at both the perfect and approximate models (16) and (18) to understand the
739 mechanism that leads to such a finite time blowup issue in the approximate model. First,
740 when u and v are at the attractor, u stays in positive values. When the amplitude of u
741 increases due to a negative value of v , both the linear and nonlinear feedback in (16) and
742 (18) will push v back to a negative value and the consequence is that v will strongly damp
743 u and decreases the amplitude of u . However, when u is negative, the nonlinear feedback
744 $-cu^2$ and the linear feedback $-cu$ will play completely different roles since $-cu^2 < 0$ while
745 $-cu > 0$. The dynamical property of the perfect dyad model (16) remains unchanged. But
746 the blowup issue appears in the approximate model (18). In fact, once u is negative, the
747 linear feedback will make v become positive. As a result, the positive anti-damping of v
748 will further increase the amplitude of u , which makes u blow up in a short time. When the
749 initial value $u(0)$ has a small amplitude (e.g., $u(0) = -0.2$), the forcing $F_u = 1 > 0$ may be
750 able to overcome the anti-damping in the short term and push the solution to the attractor.

751 But if the amplitude of $u(0)$ is large, then the role of F_u is weaker than the anti-damping
752 from v , and the solution has a much higher chance to blow up.

753 **B. the hidden variable u starting from a large positive value.**

754 Now we let the hidden variable u start from a large positive value and study how the
755 solution adjusts to the attractor. See Panels (h)–(k) in Figure 10.

756 First, with a moderately large initial condition $u(0) = 5$ as shown in Panels (h)–(i), the
757 hidden variable u using the approximate model releases to the attractor in almost the same
758 way as that using the perfect model. The trajectories of v are slightly different, but since
759 v is always very negative, the strong damping of v makes the trajectories of u in the two
760 models have very similar behavior.

761 Next, we increase the initial condition to $u(0) = 10$. Then we first notice a more significant
762 difference in the predicted trajectory of v , where in a short term $t = 0.2$ the true trajectory
763 and the perfect model prediction can reach $v = -8$ while the approximate model only allows
764 $v = -3$. This is due to the model error in the feedback terms. In fact, when u is large, $-cu^2$
765 in the perfect model will be much larger than $-cu$ in the approximate model. This leads to
766 the large error in v . As a result, the damping in the approximate model then becomes much
767 weaker compared with the perfect model. Therefore, u releases slower in the approximate
768 model (see the second row of Panel (k)). Notably, the ensemble prediction in the second
769 row of Panel (k) seems not to be too far from the truth (black dashed curve). But the truth
770 is outside 99 percentile of the prediction (the most bottom red curve) when t is between
771 $t = 0.1$ and $t = 0.5$.

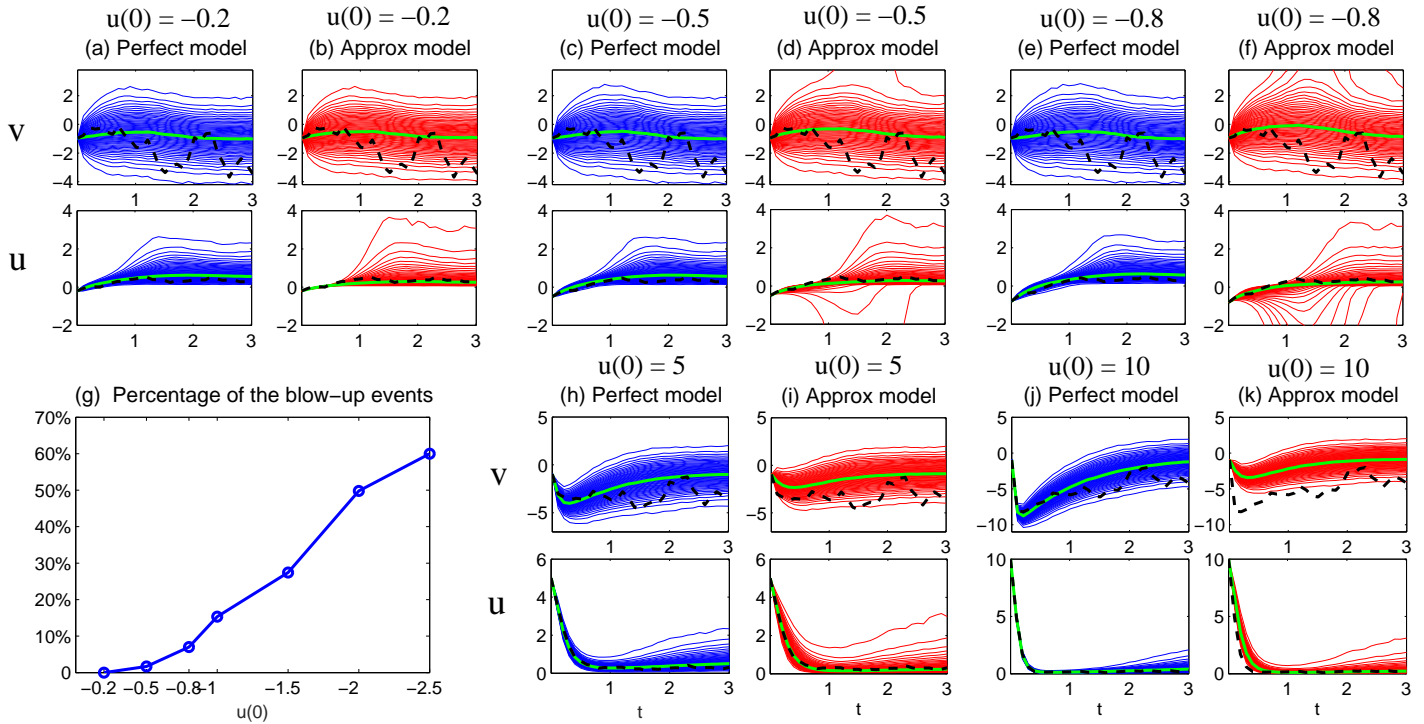


FIG. 10. Prediction of the dyad model with initial values being outside the attractor. Panels (a)–(f) and (h)–(k): Time evolution of the predicted PDFs starting from different initial phases. Each PDF is shown with 50 thin curves, which represent the 1st, 3rd, 5th, ..., 97th and 99th percentiles of the of the PDF. The green curve represents the mode of the PDF since the PDF is non-Gaussian. The black dashed curve is the true signal. In Panels (a)–(f), the hidden variable u starting from a negative value. In panels (h)–(k), the hidden variable u starting from a large positive value. Panel (g) shows the percentage of the events that blow up as a function of the initial value $u(0)$.

772 **F. Dynamical regimes with different F_u**

773 So far, we have focused on the regime with $F_u = 1$. In this subsection, the role of F_u
774 will be explored and dynamical regimes with different F_u will be studied for predicting the
775 hidden extreme events.

776 In Panel (a) of Figure 11, the trajectories of u from the perfect model (16) with different
777 F_u are shown. Here, the same random number seeds are used in generating these time series
778 for a fair comparison.

779 **Regime I:** $0.7 \leq F_u$.

780 When F_u is sufficiently large, the approximate model with the linear feedback (18) is a
781 suitable model for predicting the hidden extreme events.

782 **Regime II:** $0 \leq F_u < 0.3$.

783 When F_u approaches zero, the intermittent events in u can have both signs. As was
784 discussed in Section V E, when u is negative, the linear feedback $-cu$ in (18) will play a
785 significant different role compared with the nonlinear feedback $-cu^2$ in the perfect model
786 (16). In fact, the linear feedback $-cu$ becomes positive and make v to be positive. Then
787 the anti-damping of v in the process of u leads to the finite time blowup. Therefore, we
788 conclude that using the approximate model (18) with a linear feedback to predict the hidden
789 extreme events in u requires that the forcing F_u in the perfect dyad model cannot be too
790 small. If the forcing F_u in the perfect dyad model is too small, then the approximate model
791 does not have a mechanism to recover the intermittent events in u when u is negative. A
792 new approximate model that has skill in capturing the extreme events with both signs needs
793 to be developed.

794 **Regime III:** $0.3 \leq F_u < 0.5$.

795 When $F_u \geq 0.3$, the intermittent events in the true signal of u only occur in the positive
796 phase. However, the true trajectory of u still goes below 0 quite frequently (with small
797 amplitudes). Panel (c) of Figure 11 shows the data assimilation of u using the approximate
798 model with the linear feedback (18), where the parameters are re-estimated based on the
799 observed signal of v in $F_u = 0.3$ regime. One important result is that the assimilated state
800 of u can occasionally become quite negative! In fact, as is shown in Panels (b)–(c), before
801 the assimilated u goes to a negative value, the signal of v is large and positive while u is

802 nearly zero. Therefore, when the trajectory of u becomes slightly negative in the true signal,
803 the anti-damping v will amplify the negative phase of u . Since the positive forcing $F_u = 0.3$
804 here is pretty weak, this forcing is unable to push u back to the attractor with positive
805 values immediately and therefore the assimilated state of u will stay in the negative phase
806 for a while. According to the discussions in Section V E, if the prediction starts with a large
807 negative value of u , then even for a short term, the prediction using the approximate model
808 may suffer from a short-term blowup¹¹⁰.

809 **Regime IV:** $0.5 \leq F_u < 0.7$.

810 Now the data assimilation results using the approximate model (18) provides the state
811 of u that is always positive. Thus, there will be no issue in data assimilation. However, as
812 shown in Panels (d)–(f) of Figure 11, the approximate model can still suffer from a long (but
813 finite) time blow up issue. This is again related to the insufficient strength of F_u . In fact,
814 the trajectory of u still has some chances to become slightly negative and the corresponding
815 values of v at these time instants are usually large. Therefore, the anti-damping v and the
816 forcing F_u in the process of u compete with each other. If the strength of forcing is not
817 strong enough, then for some events, the anti-damping can results in the blowup issue of
818 u ¹¹⁰.

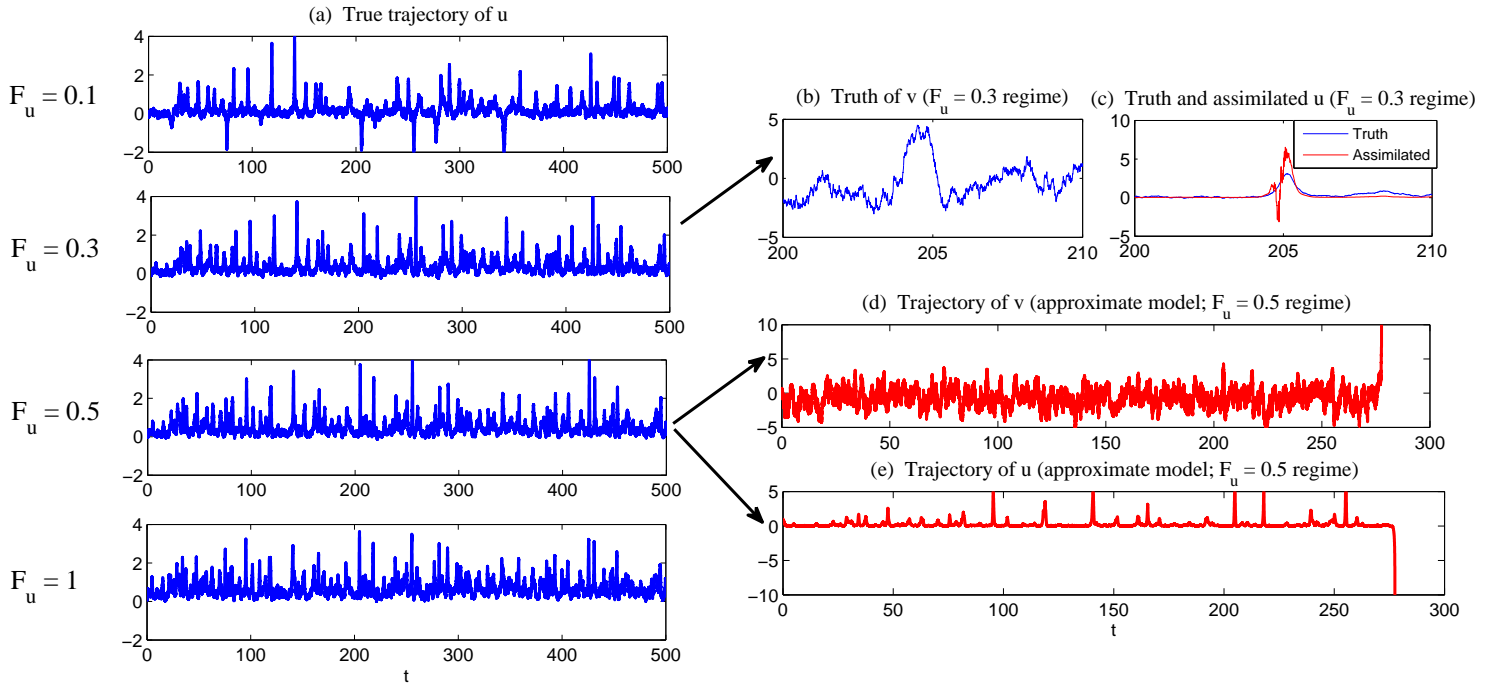


FIG. 11. Dynamical regimes with different F_u . Panel (a): trajectories of u from the perfect model (16) with different F_u are shown. Here, the same random number seeds are used in generating these time series for a fair comparison. Panel (b): True signal of v in $F_u = 0.3$ regime from the perfect model. Panel (d): True signal of u in $F_u = 0.3$ regime from the perfect model (blue) and the assimilated posterior mean using the approximate model. Panels (e) and (f): trajectories of the approximate model with the estimated parameters from the observed true signal of v in $F_u = 0.5$ regime.

819 VI. THE LORENZ 63 MODEL

820 In many applications with chaotic or turbulent phenomena, due to the incomplete knowl-
821 edge of the underlying dynamics, noise inflation is often incorporated into the dynamical
822 processes^{21,104,111}. The enhanced noise plays the role of parameterizing small-scale fluctua-
823 tions, which helps increase the variability of the system and has a potential of improving
824 the data assimilation and prediction skill. Yet, it has not been well understood the effect of
825 noise inflation in the extreme events prediction. Therefore, in this and the next two sections
826 (Section VII and Section VIII), noise inflation will be incorporated into the dynamical sys-
827 tems for testing the ensemble prediction skill of the observed and hidden extreme events as
828 well as other non-Gaussian characteristics. The difference between the studies in these three
829 sections is as follows. In Section VII and Section VIII, the noise inflation will be combined
830 with various effective and practical strategies for developing effective and simple approxi-
831 mate models for improving the prediction of the extreme events resulting from complicated
832 turbulent dynamical systems with regime switching. In this section, the chaotic Lorenz 63
833 model is used as a testbed to understand the skill of the extreme events predictions, where
834 the inflated noise acts as the only source of the model error.

835 A. The perfect and approximate models

836 The model considered in this section is the Lorenz 63 model⁷². It is a simplified mathe-
837 matical model for atmospheric convection with chaotic behavior. The equations relate the
838 properties of a two-dimensional fluid layer uniformly warmed from below and cooled from
839 above. In particular, the equations describe the rate of change of three quantities with
840 respect to time: x is proportional to the rate of convection, y to the horizontal tempera-
841 ture variation, and z to the vertical temperature variation. The constants σ , ρ , and β are
842 system parameters proportional to the Prandtl number, Rayleigh number, and certain phys-
843 ical dimensions of the layer itself¹¹². The Lorenz 63 model is also widely used as simplified
844 models for lasers, dynamos, thermosyphons, electric circuits, chemical reactions and forward
845 osmosis^{113–119}.

Here, we study a slightly different version of the original Lorenz 63 model by adding a small noise into the x process. We also assume that only a short trajectory of x is observed

as the training data while y and z are the unobserved variables. The model reads:

$$dx = \sigma(y - x)dt + \sigma_x dW_x, \quad (21a)$$

$$dy = (x(\rho - z) - y)dt, \quad (21b)$$

$$dz = (xy - \beta z)dt, \quad (21c)$$

846 The small noise here can be regarded as the observational or measurement uncertainty. It
 847 also helps prevent the singularity in the data assimilation formula in (3), which requires a
 848 non-zero noise in the observational process. Nevertheless, with a small noise coefficient, the
 849 dynamical behavior of the model in (21) remains almost the same as the original noise-free
 850 Lorenz 63 model. Below, we always take $\sigma_x = 1$, which is a sufficiently small value. The
 851 other parameters that are used to generate the true signals of (21) are

$$\sigma = 10, \quad \rho = 28, \quad \beta = 8/3. \quad (22)$$

852 These are the classical choices of the Lorenz 63 model. Figure 12 shows the trajectories,
 853 PDFs and phase plots of the Lorenz 63 model (21), where the butterfly profile in the phase
 854 plots and the chaotic features in the model trajectories are clearly demonstrated. Notably,
 855 there are quite a few extreme events that appear in all the three components due to the
 856 fact that one of the Lyapunov exponents of the Lorenz 63 system is positive. These extreme
 857 events occur when the system states switch between the two branches of the “butterfly
 858 wings”.

859 The short trajectory of x in Panel (a) of Figure 12 with only 50 units will be used as the
 860 observed training data for the approximate models below.

861 **The approximate models.**

862 Below, we aim at understanding the model error that comes from the noise inflation. To
 863 this end, it is natural to propose the following approximate model,

$$\begin{aligned} dx &= \sigma(y - x)dt + \sigma_x dW_x, \\ dy &= (x(\rho - z) - y)dt + \sigma_y dW_y, \\ dz &= (xy - \beta z)dt + \sigma_z dW_z, \end{aligned} \quad (23)$$

864 where the noise coefficients σ_y and σ_z are given and fixed empirically, which account for the
 865 noise inflation in the hidden variables. Note that here the deterministic parts in the perfect
 866 model (21) and the approximate model (23) are the same, which is not always the case in

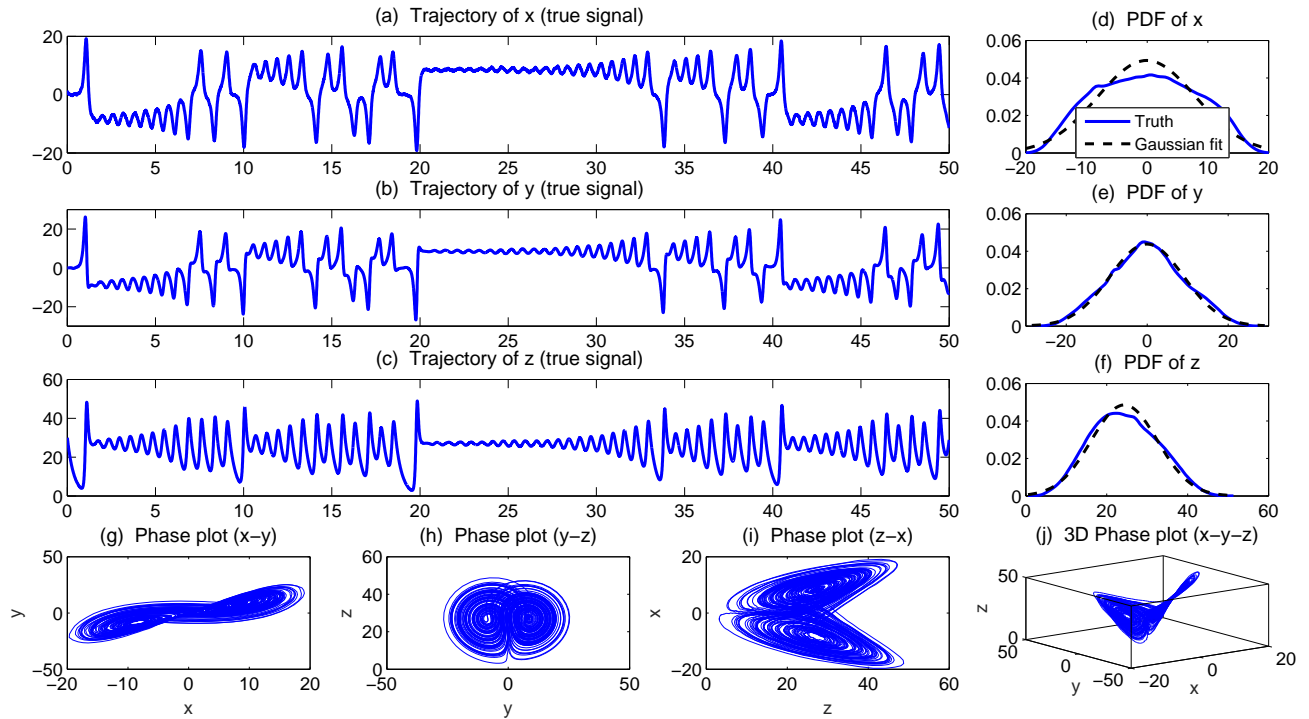


FIG. 12. Trajectories (Panels (a)–(c)), PDFs (Panels (d)–(f)) and the phase plots (Panels (g)–(j)) of the noisy Lorenz 63 model (21) with parameters in (22) and a small noise coefficient $\sigma_x = 1$.

867 real applications where noise inflation is often used to compensate other model errors and
 868 simplifications. Nevertheless, the setup here allows us to understand the model error that
 869 comes purely from the noise inflation and its effect on the prediction skill.

870 The noise coefficient σ_x will be estimated from the parameter estimation algorithm. Note
 871 that since σ_x is associated with the quadratic variation of the continuously observed training
 872 data, a prescribed value with inflation may lead to pathological behavior of the parameter
 873 estimation. Depending on the level of noise inflation, we consider the following three ap-
 874 proximate models,

Approximate model I.	Small noise inflation:	$\sigma_y = \sigma_z = 1,$	
Approximate model II.	Moderate noise inflation:	$\sigma_y = \sigma_z = 3,$	(24)
Approximate model III.	Large noise inflation:	$\sigma_y = \sigma_z = 5.$	

875 **B. Parameter estimation**

876 In the approximate models, there are four parameters to be estimated: ρ, σ, β and σ_x .
 877 Here the parameter estimation algorithm as described in Section III D is run up to $K = 15000$
 878 steps and the averaged values of the trace plots from $k = 5000$ to $k = 15000$ are used as the
 879 estimated parameters, which are:

$$\begin{aligned}
 \text{Approx model I:} \quad & \rho = 27.48, \quad \sigma = 10.34, \quad \beta = 2.70, \quad \sigma_x = 1.03, \\
 \text{Approx model II:} \quad & \rho = 31.04, \quad \sigma = 9.051, \quad \beta = 2.33, \quad \sigma_x = 1.06, \quad (25) \\
 \text{Approx model III:} \quad & \rho = 34.17, \quad \sigma = 7.525, \quad \beta = 2.20, \quad \sigma_x = 1.08.
 \end{aligned}$$

880 Note that due to the model error from noise inflation, the estimated parameters in the
 881 approximate models are not exactly the same as those in the perfect model. In particular,
 882 with the increase of the noise coefficients σ_y and σ_z , the estimated parameter ρ and σ seem
 883 to be more different compared with the one in the perfect model in order to compensate the
 884 model error.

885 **C. Data assimilation**

886 Figure 13 shows the data assimilation results using the approximate model (23) with the
 887 estimated parameters, where the true signal of the observed variable x is generated using
 888 the perfect model (21).

889 In the approximate model I, due to the small model error in the inflated noise coefficients,
 890 the assimilated values of y and z are nearly perfect and the uncertainty reflected by the
 891 posterior variance in both variables is small. In the approximate model II, the assimilated
 892 values of y are still quite accurate but those of z show some errors where the mean state of
 893 z has a slight shift towards the positive value. Such a bias in the assimilated posterior mean
 894 state is possibly due to the fact that the noise σ_y leads to the change of the mean value of
 895 xy in z process since x and y are highly correlated. On the other hand, x and z are not so
 896 closely correlated, and therefore the mean value of xz that contributes to the mean state of
 897 y is hardly polluted by the noise. Finally, in the approximate model III, where the inflated
 898 noise coefficients are large, there are some non-negligible errors in the assimilated states of
 899 z and the associated uncertainty increases as well. Nevertheless, despite such a mean state
 900 shift, the overall patterns and amplitudes of z are assimilated quite well.

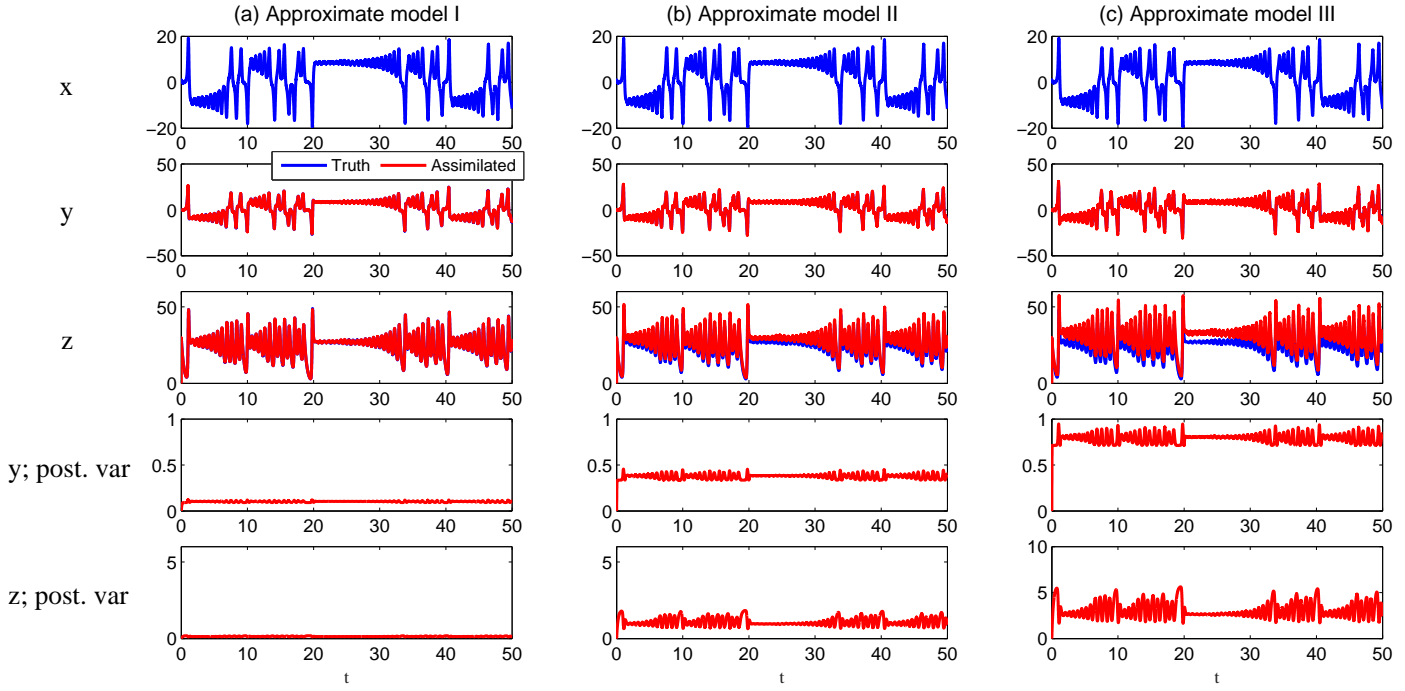


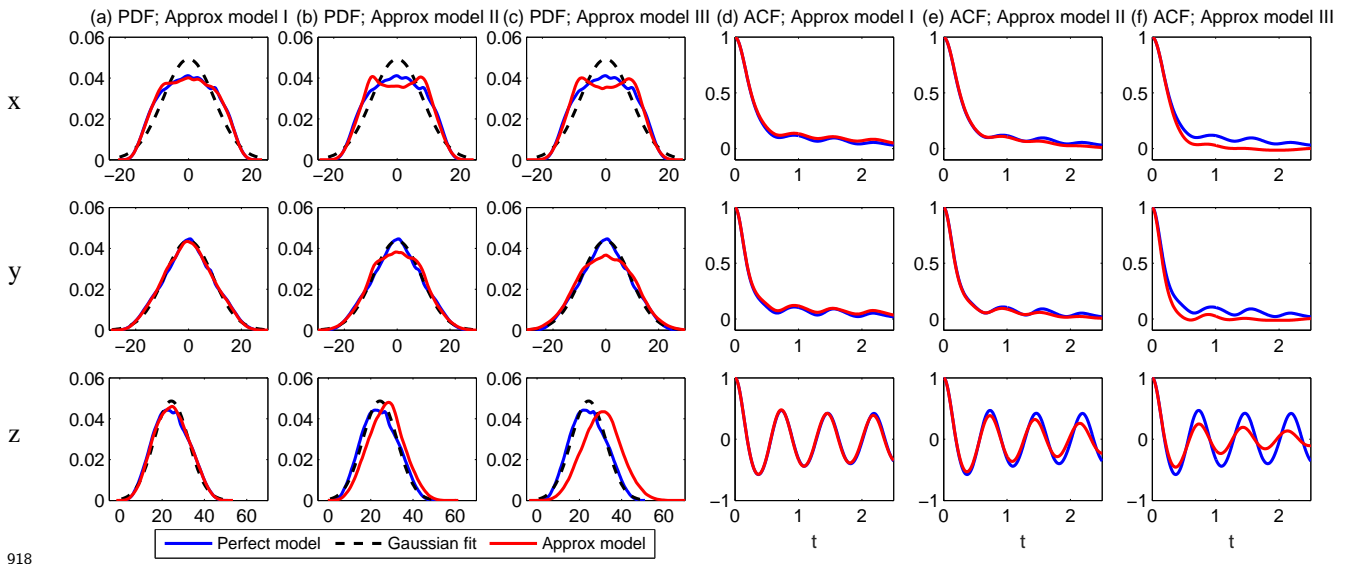
FIG. 13. Data assimilation using the approximate model (23) with different noise inflation levels (24). The first row shows the true trajectory of x . The second and third rows show the true signals of y and z as well as the posterior mean estimations from data assimilation (red). The fourth and fifth rows show the posterior variance of y and z , respectively.

902 **D. Long-range forecast**

903 To quantify the long-range forecast skill, the comparison of the equilibrium PDFs and
 904 the ACFs between the perfect model and approximate models is shown in Figure 14.

905 First, all the three approximate models are able to capture the equilibrium non-Gaussian
 906 PDFs of both x and y with high accuracy, where the information model error in the equilib-
 907 rium PDF $\mathcal{P}(p_{eq}, p_{eq}^M) \leq 0.05$ is tiny even using the approximate model III. For the variable
 908 z using the approximate model III, the error is slightly larger $\mathcal{P}(p_{eq}, p_{eq}^M) = 0.28$ but is still
 909 acceptable. Such a model error is due to the fact that the PDF associated with z using the
 910 approximate model has a mean shift compared with the truth, which has already been seen
 911 in the data assimilation results.

912 Next, the approximate models and the perfect model overall share quite similar ACFs,
 913 indicating similar time evolution behavior (at least up to the second order statistics in time).
 914 In particular, the direct relaxation of the ACFs of x, y and the oscillated relaxation of that of
 915 z are both captured by the approximate models. The only non-negligible difference appears
 916 in the ACF of z when the noise inflation level is large, i.e., in approximate model III, where
 917 the approximate model has a slightly faster decaying ACF.



918
 919 FIG. 14. Comparison of the PDFs and the ACFs of the perfect model and the approximate model.

920 E. Short- and medium-range forecasts

921 To study the short- and medium-range forecasts, we show in Figure 15 three skill scores
922 of the predictions as a function of lead time. Two of them, namely the RMSE and the Corr,
923 are the classical path-wise measurements while the third one is an information criterion, that
924 is, the relative entropy (7) between the PDF of the predicted time series and that of the
925 truth. In order to distinguish the errors due to the noise inflation and the initial uncertainty
926 with data assimilation, we show the predictions using the approximate model with either
927 assimilated initial conditions (ICs) or with perfect initial conditions. All the predictions
928 here are based on the ensemble mean, which is the average of 50 ensemble members.

929 Columns (a)–(b) and Columns (d)–(e) show the RMSE and Corr of the predictions using
930 the approximate models I and II, respectively. These path-wise measurements indicate that
931 the skillful predictions of the approximate models regardless of using perfect or assimilated
932 initial conditions are up to nearly 3 time units. However, the conclusion based on these
933 path-wise measurements can be misleading in this strongly chaotic system. In Columns
934 (g)–(h), the relative entropy has a significant increase as the lead time, especially using the
935 approximate model II. This implies certain non-negligible errors are not captured by the
936 two path-wise measurements. To see such errors, the ensemble mean prediction using the
937 approximate models (green) and the truth (blue) at lead time $t = 1$ are compared in Figure
938 16. Both the trajectories and the PDFs are shown in order to compare the path-wise and the
939 information measurements. Note that only the Gaussian fits of the PDFs are shown here for
940 the purpose of comparing the variance in the truth and the predicted PDFs which reflects
941 the skill of capturing the amplitudes especially those of the extreme events. In Column
942 (b) of Figure 16, it is shown that although the patterns of the predicted signal are quite
943 consistent with the truth, the amplitudes of all the extreme events are underestimated to
944 some extent. Thus, the predicted PDF has a narrower shape compared with the truth. Such
945 a phenomenon becomes more significant in Column (c) of Figure 16 where the approximate
946 model III is used. At lead time $t = 1$, despite that $\text{Corr} \approx 0.8$ for x and y and $\text{Corr} \approx 0.5$
947 for z remain skillful, the large values of the relative entropy clearly indicate the discrepancy
948 between the predicted PDF and the truth, which is due to the fact that the amplitudes
949 of the extreme events are severely underestimated. These facts conclude the importance of
950 using the information criterion in quantifying the model error in the PDFs in addition to

951 the path-wise measurements.

952 Figure 17 includes a case study of the ensemble prediction for short- and medium range
953 forecasts using approximate models I and III starting from $t = 16.5$. First, the ensemble
954 mean (green) using the approximate model I is skillful up to 2.5 units of the lead time
955 while that using the approximate model III has a much shorter skillful prediction. These
956 are consistent with the results shown in Figure 16. Next, the uncertainty of the prediction
957 is reflected in the ensemble spread. It is clear that in the perfect model prediction the
958 ensembles do not spread out until $t = 19$ while those in the approximate models start
959 spreading out around $t = 17$. This is obviously due to the fact that the noise level is
960 higher in the approximate models. Using approximate model I, despite some members
961 give false prediction due to the intrinsic chaotic behavior, most of the ensemble members
962 are still able to follow the true trajectories, which also results in the skillful ensemble mean
963 prediction. However, using the approximate model III, both the large noise inflation and the
964 initial uncertainty due to the data assimilation lead to a quick divergence of the ensembles.
965 The ensemble spread is able to tell the uncertainty but the ensembles reach the attractor
966 much faster than those using the perfect model and therefore the ensemble mean using the
967 approximate model losses its skill.

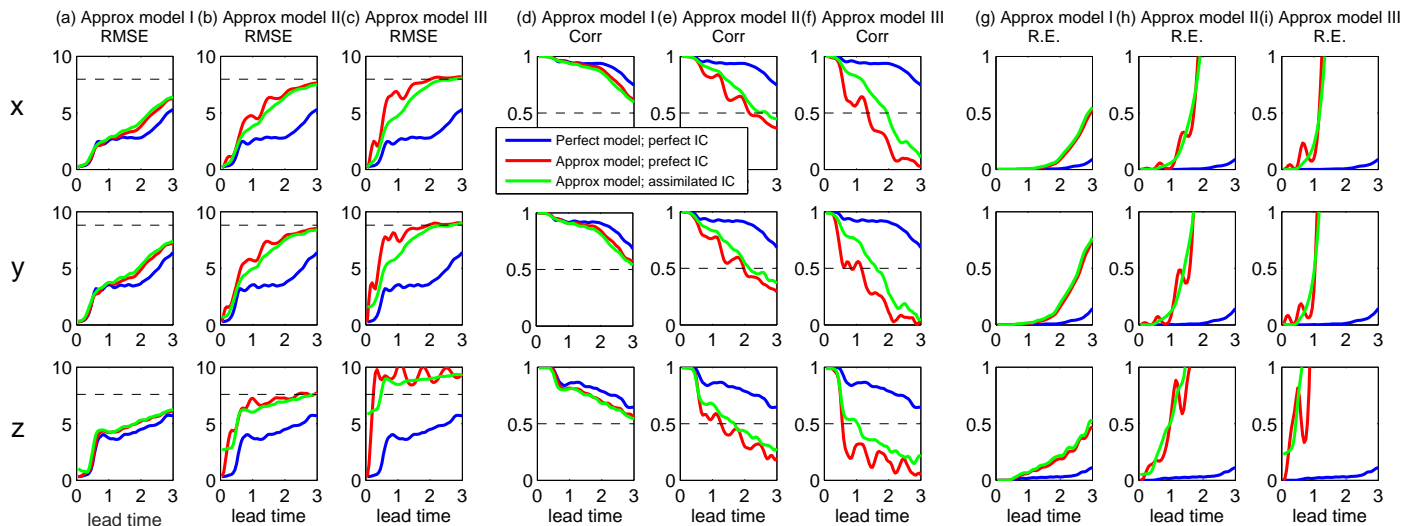


FIG. 15. RMSE (Panels (a)–(c)), Corr (Panels (d)–(f)) and relative entropy (R.E.; Panels (g)–(i)) as a function of lead time for short- and medium-range forecasts using the perfect model (21) (blue) and the three approximate models (23)–(24) with perfect initial conditions (red) and assimilated initial conditions (green). The prediction here is based on the ensemble mean. The dashed black lines in the RMSE panels show one standard deviation of the true signal and those in the Corr panels show the $\text{Corr} = 0.5$ threshold.

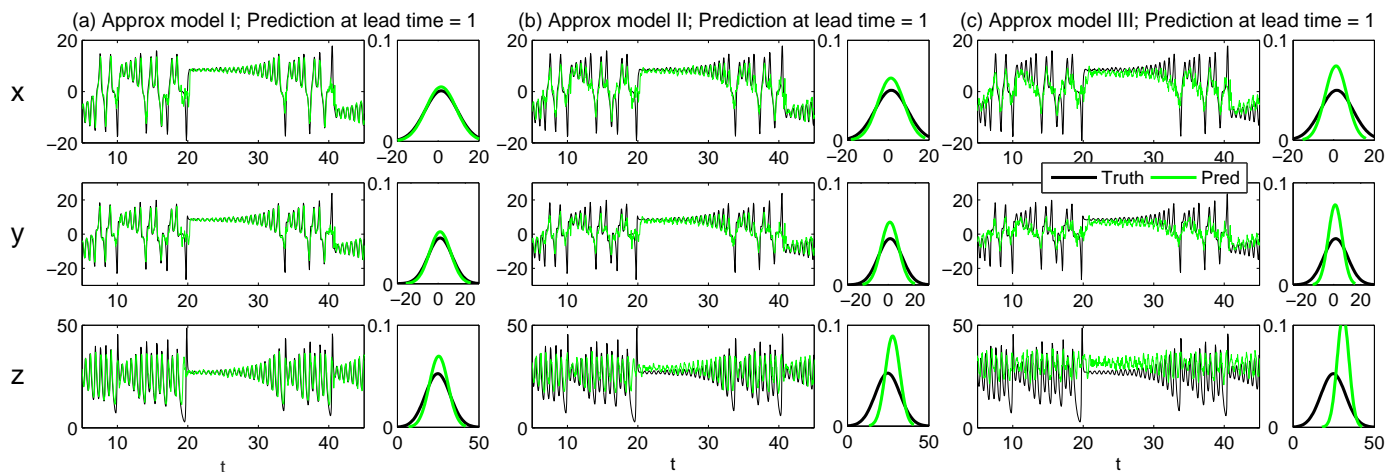


FIG. 16. Comparison of the ensemble mean prediction using the approximate models and the assimilated initial conditions (green) with the truth (blue) at lead time $t = 1$. In each panel, both the trajectories and the Gaussian fits of the PDFs are shown.

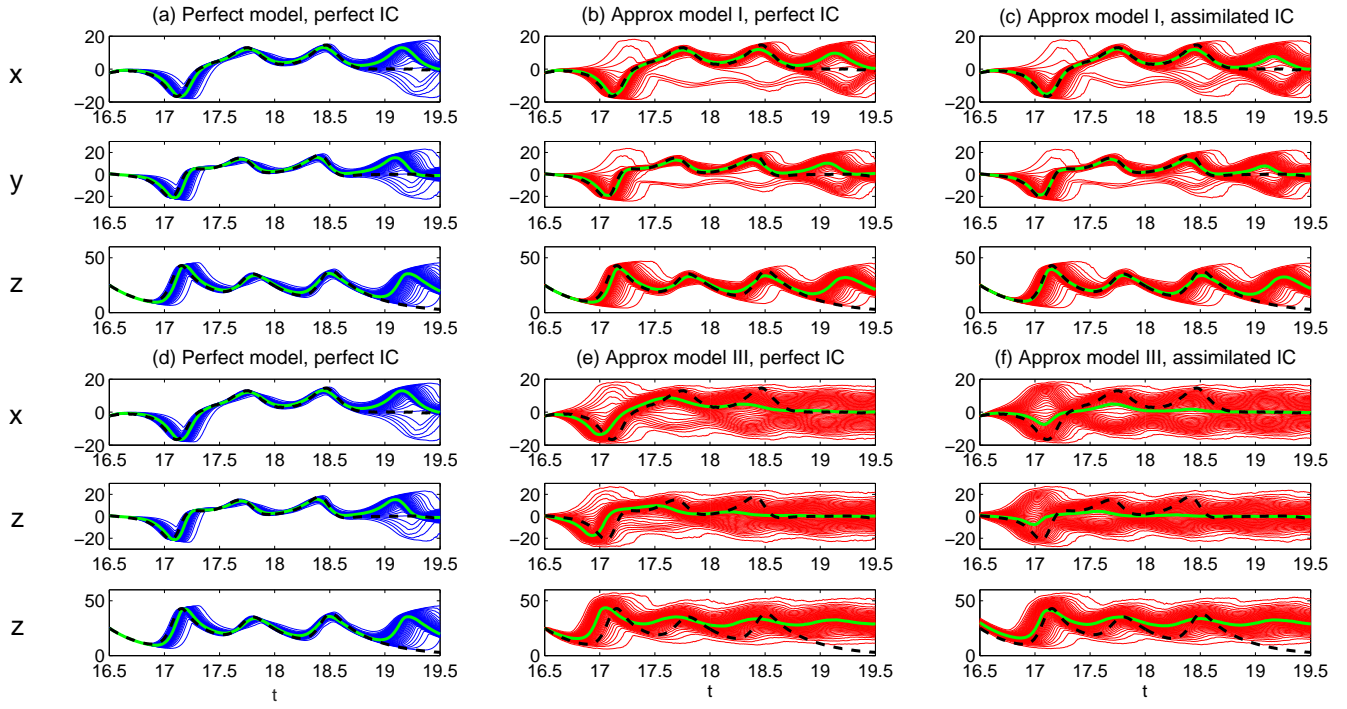


FIG. 17. Case studies of the ensemble forecasts. Panels (a)–(c): a case study using the approximate model I. Panels (d)–(f): the same case study using the approximate model III. Each subpanel shows the time evolution of the prediction, as represented by the time dependent PDF of the ensemble forecast. Note that each PDF is shown with 50 thin curves, which represent the 1st, 3rd, 5th, . . . , 97th and 99th percentiles of the of the PDF. The green curve represents the mode of the PDF since the PDF is non-Gaussian. The black dashed curve is the true signal.

969 **F. Prediction with an initial value starting outside the attractor**

970 Panels (a)–(f) and Panels (g)–(l) in Figure 18 show the prediction where the initial values
971 of the observed variables $x(0) = 150$ and those of the unobserved ones $y(0) = z(0) = 150$
972 are outside the attractor, respectively.

973 The skill of capturing the transition behavior of the approximate models depends on the
974 model error in the noise inflation. The approximate model I behaves almost the same as the
975 perfect model due to its small noise inflation. The approximate model II is able to capture
976 the transition behavior in short and medium ranges if the initial values of y and z are off
977 the attractor. However, it fails to predict the two hidden variables after a very short period
978 if the initial value of x is off the attractor. On the other hand, the approximate model
979 III, which has the largest noise inflation coefficients, only has skillful prediction for a short
980 period no matter which variable starts from a value that is outside the attractor.

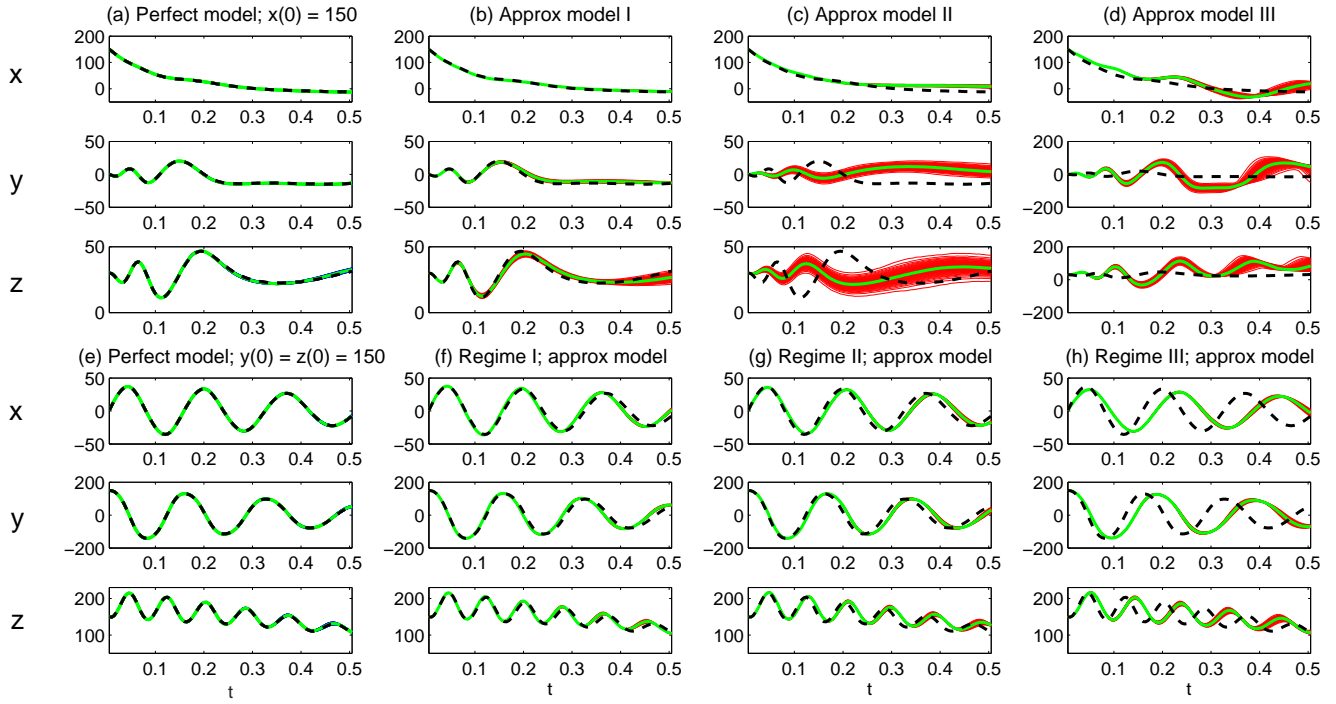


FIG. 18. Prediction with an initial value starting outside the attractor. Panels (a)–(d): Prediction where x starts at $x = 150$, which is a value that is off the attractor. Panels (e)–(h): Prediction where both y and z start at $y = z = 150$, which are values that are off the attractor. Each PDF is shown with 50 thin curves (blue for the perfect model and red for the approximate model), which represent the 1st, 3rd, 5th, \dots , 97th and 99th percentiles of the of the PDF. The green curve represents the mode of the PDF since the PDF is non-Gaussian. The black dashed curve is the true signal.

981 **VII. A PARADIGM MODELS FOR TOPOGRAPHIC MEAN FLOW**
 982 **INTERACTION WITH REGIME SWITCHING BEHAVIOR**

983 Regime switching between multiple metastable states is a key feature in many nonlinear
 984 turbulent dynamical systems^{120–122}. One example is the atmospheric flow regimes, which rep-
 985 resent the recurrence of certain flow structures despite the intrinsic chaotic behavior of the
 986 underlying system. The existence of persistent or recurrent weather patterns¹²³ with block-
 987 ings is one of the most pronounced illustrations of synoptic-scale circulation regimes^{124,125}
 988 while different circulation regimes and their switching were also found in planetary-scale
 989 patterns^{126,127}. The metastable states have their unique dynamical behavior and the regime
 990 switching often triggers extreme events and other important nonlinear phenomena. Notably,
 991 the regimes can appear even though the observed data have a nearly Gaussian probability
 992 distribution^{122,128,129}. Due to the highly complex nature of these regimes and their switching
 993 behavior as well as only the availability of the partial observations, it is important to develop
 994 suitable approximate models for capturing both the dynamical and statistical features of the
 995 regime switching and for predicting the associated extreme events. In this section, we con-
 996 centrate on the development of nonlinear low-order models to achieve the above tasks, where
 997 the topographic effect is regarded as the result of random structures from either atmosphere
 998 or ocean in intermediate and small scale.

999 **A. The perfect model**

Consider the barotropic quasi-geostrophic equations²,

$$\frac{\partial q}{\partial t} + \nabla^\perp \psi \cdot \nabla q + u(t) \frac{\partial q}{\partial x} + \beta \frac{\partial \psi}{\partial x} = 0, \quad (26a)$$

$$q = \Delta \psi + h, \quad (26b)$$

$$\frac{du}{dt} = \int h \frac{\partial \psi}{\partial x}. \quad (26c)$$

1000 This is an ideal model to study the complex nonlinear interaction of the large-scale and
 1001 the small-scale flow and the role of the topography. The model exhibits a regime switching
 1002 behavior with blocked and unblocked zonal flow structure despite that the associated PDF
 1003 of the zonal flow has only a single modal. The study of this model for understanding its
 1004 mathematical properties, developing reduced order models and uncertainty quantification

1005 can be found in a series of papers^{2,36,53,130,131}. In particular, rigorous statistical bounds in
 1006 quantifying the uncertainty for the ensemble prediction of barotropic flow over topography
 1007 has been shown in a recent paper¹³².

1008 In this model, the small-scale flow is given in terms of the stream function ψ , and q is the
 1009 small-scale potential vorticity. The large-scale velocity field only has the zonal component
 1010 $u(t)$, and the topography is given by the function $h = h(x, y)$. The parameter $\beta > 0$ is the
 1011 contribution from the beta-plane effect. Both the small-scale potential vorticity q and the
 1012 small-scale stream function ψ , as well as the topography h , are assumed to be 2π -periodic
 1013 functions in both variables x and y with zero average. The large-scale velocity $u(t)$ is strongly
 1014 coupled with the small-scale flow through equation (26c), where the bar across the integral
 1015 sign indicates that the integral has been normalized by the area of the domain of integration.

1016 Below, we consider a special situation to the full nonlinear system, which inherits the
 1017 nonlinear coupling of the small-scale flow with the large-scale mean flow via topographic
 1018 stress. The model is named as the *layered topographic equations*. Here the topography is
 1019 layered in the fixed direction $\vec{l} = (l_x, l_y)$. We assume that both ψ and q only depend on
 1020 $\xi = \vec{l} \cdot \vec{x}$ with $\vec{x} = (x, y)$. One key feature of the layered topographic equations is that
 1021 the small-scale nonlinear term in (26a), $\nabla\psi \cdot \nabla^\perp q$, is identically zero. Nevertheless, the
 1022 nonlinear coupling due to topographic stress remains and is responsible for much of the
 1023 complex behavior. Without loss of generality we can always rescale the system with $l_x \neq 0$
 1024 to align to a special case with $\vec{l} = (1, 0)$.

1025 In such a situation, the Fourier expansion of ψ and h are given by

$$\begin{aligned} \psi(x, y, t) &= \sum_{k \neq 0} \psi_k(t) e^{ik\vec{l} \cdot \vec{x}}, \\ h(x, y) &= \sum_{k \neq 0} h_k e^{ik\vec{l} \cdot \vec{x}}, \end{aligned} \tag{27}$$

1026 where we have assumed that the topography has zero mean with respect to spatial average,
 1027 that is $h_0 = 0$. Substituting the ansatz (27) into (26) and adding stochastic forcing and
 1028 damping, we arrive at the layered topographic equations in Fourier form,

$$\begin{aligned} \frac{d\psi_k}{dt} &= -d_k \psi_k + ikl_x \left(\frac{\beta}{k^2 |\vec{l}|^2} - u \right) \psi_k + i \frac{kl_x}{k^2 |\vec{l}|^2} h_k u + \sigma_k \dot{W}_k, \\ \frac{du}{dt} &= -d_u u - il_x \sum_{k \neq 0} kh_k \psi_k^* + \sigma_u \dot{W}_u, \end{aligned} \tag{28}$$

1029 where $*$ denotes the complex conjugate. In (28), $\psi_k, k = 1, 2, \dots, \Lambda$ are the stream functions
 1030 and u is the large-scale zonal velocity.

1031 In the study here, we adopt $\Lambda = 10$ and therefore in total there are 21 modes in the
 1032 model (28), where 1 mode u represents the large-scale zonal flow. The other 20 modes
 1033 are for the small-scale stream functions with $k = \pm 1, \dots, \pm 10$, which based on the layered
 1034 topographic functions determine the meridional flows. We assign the following function for
 1035 the topography,

$$h(x) = H_1 \left(\cos(x) + \sin(x) \right) + H_2 \left(\cos(2x) + \sin(2x) \right) - \frac{i}{2} \sum_{3 \leq k \leq \Lambda} \frac{e^{i(kx + \theta_k)}}{k^p} + c.c., \quad (29)$$

1036 where H_1 and H_2 are associated with the leading two Fourier modes $k = \pm 1, \pm 2$ while
 1037 the remaining part in (29) represents the amplitudes of the topography for other Fourier
 1038 modes. Here θ_k are random phase and p is a power that controls the effects of the small-scale
 1039 topography. The topography plays an important role in altering the stream functions. With
 1040 a simple manipulation, it is easy to show that the topographic functions associated with the
 1041 first two Fourier modes are

$$h_1 = H_1/2 - H_1/2i, \quad \text{and} \quad h_2 = H_2/2 - H_2/2i. \quad (30)$$

1042 The other h_k can also be written explicitly using (29). The following parameters are adopted
 1043 in the study here. The beta-plane effect is $\beta = 2$. The coefficients of the topography are
 1044 $H_1 = 1$ and $H_2 = 1/2$. The damping coefficients are chosen as

$$d_k = d_u = 0.0125, \quad (31)$$

1045 which represents a time scale of relaxation time of roughly 80 time units. Such a choice allows
 1046 a relatively slow (but not infinitely slow) mixing of the system. With different choices of
 1047 the stochastic noise, the system can also have fast mixing rate. Finally, the stochastic noise
 1048 coefficients are chosen as follows,

$$\sigma_u = \sigma_1 = \sigma_2 = \frac{1}{20\sqrt{2}}, \quad \sigma_k = \frac{1}{20\sqrt{2}} \frac{1}{k^p}, \quad \text{for } p = 3, \dots, \Lambda. \quad (32)$$

1049 **Dynamical regimes.**

1050 Two dynamical regimes will be studied below, which correspond to different values of p
 1051 with $p = 1$ and $p = 0.5$. Note that the ‘‘dynamical regimes’’ here should not be confused

1052 with the “regime switching”. The latter stands for the switching of the model variables
 1053 between different values or states within a given dynamical regime.

1054 Figure 19 shows the time series of the zonal flow u , its associated PDFs and ACFs as
 1055 well as the accumulated energy in the small-scale stream functions. Here, the accumulated
 1056 energy $E[\psi_{1:s}]$ is defined as

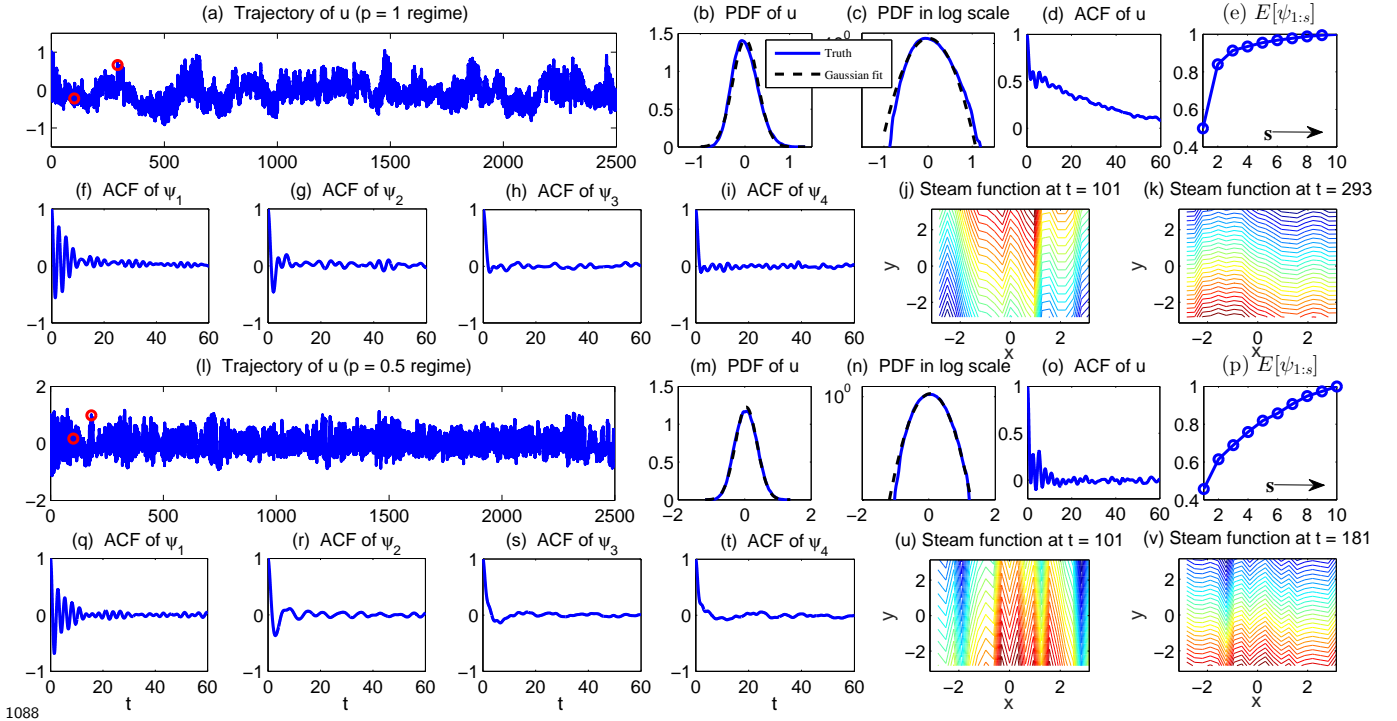
$$E[\psi_{1:s}] = \sum_{k=1}^s |\psi_k|^2. \quad (33)$$

1057 In the regime with $p = 1$, the trajectory of the zonal velocity u switches between roughly
 1058 two different states and it stays in each state for a while before switching to the other
 1059 (see Panel (a)). One state with positive u corresponding to the eastward zonal flow contain
 1060 extreme events. Despite the nearly “two-state” time series, the associated PDF of u is single
 1061 modal and is slightly skewed where a single (positive) side fat tail correspond to extreme
 1062 events for the eastward zonal flow. Note that the regime (state) switching behavior with
 1063 such a single modal distribution has been systematically studied in¹²². Despite the single
 1064 modal distribution, the ACF is highly different from a Gaussian model with exponential
 1065 decay. In fact, the ACF here first experiences a sharp decrease to $\text{ACF} = 0.5$ and then it
 1066 decays slowly with almost a linear decaying rate to zero. The total decaying time is about
 1067 60 time units. On the other hand, regarding the small-scale stream functions ψ_k , the leading
 1068 two modes contain about 84% of the total energy. The ACF associated with ψ_1 has a strong
 1069 oscillation with a long memory while that associated with ψ_2 only has a weak oscillation.
 1070 For modes ψ_k with $k \geq 3$, the ACFs decay quite fast.

1071 Next, in $p = 0.5$ regime, the trajectory of the zonal velocity u has a relatively strong
 1072 mixing rate. The direction of the zonal velocity alternates between eastward and westward
 1073 quite frequently. Despite the Gaussian statistics, the dynamical regime is still chaotic. The
 1074 ACF associated with u now behaves in a very different way, where it oscillates and decays
 1075 quickly to zero. The leading two modes of the small-scale stream functions ψ_1 and ψ_2 contain
 1076 about 61% of the total energy, and the ACFs associated with ψ_k with $k \geq 3$ now decay more
 1077 slowly compared with those in $p = 1$ regime.

1078 Notably, in both regimes, the total flow field alternatives between zonally blocked and
 1079 unblocked patterns as shown in Panels (j)–(k) and (u)–(v). Recall in (27) the topographic
 1080 effect is imposed on the layered modes with $\vec{l} = (1, 0)$. This implies that the contributions
 1081 of all the small-scale stream functions ψ_k are on the meridional flows while the zonal flow

1082 is driven by the large-scale zonal mode u . As a consequence, when the total flow field is
 1083 zonally blocked, the large scale zonal velocity $u = 0$ and the total energy lies in the small-
 1084 scale stream functions (see Panels (j) and (u)). Similarly, when the zonal flow becomes
 1085 dominant, its kinetic energy accounts for a large portion of the total energy (see Panels (k)
 1086 and (v)). Therefore, the regime switching not only alters the flow patterns but also adjusts
 1087 the energy contributions in the total flow field.



1089 FIG. 19. Dynamical regimes of the layered topographic model (28). Panels (a)–(k): regime with
 1090 $p = 1$. Panels (a)–(d) show the time series, PDF, PDF in logarithm scale and the ACF of u . Panel
 1091 (e) shows the accumulated energy $E[\psi_{1:s}]$ defined in (33). Panels (f)–(i) show the ACFs of the
 1092 first four stream functions. Panels (j)–(k) show the total streamline at two different time instants
 1093 marked in red dots in Panel (a). At these two time instants, the model shows the blocked and
 1094 unblocked zonal flow structure, respectively. Panels (l)–(v) are similar to (a)–(k) but for regime
 1095 with $p = 0.5$.

1096 **B. The approximate model**

1097 Recall that the perfect model (28) has a 21-degree of freedom. The approximate model
 1098 developed here has a much simpler form, which includes only the zonal flow u and the leading
 1099 two Fourier wavenumbers (with $k = \pm 1$ and ± 2).

1100 The motivation of such a choice of the approximate model comes from the rapid decay
 1101 of the ACFs associated with the small-scale stream functions. In fact, as shown in Figure
 1102 19, the stream functions ψ_k with $k = 3, \dots, 10$ decorrelate very fast while ψ_1 has a much
 1103 longer relaxation time and ψ_2 also has some memory. Therefore, it is natural to retain the
 1104 dynamics of the leading two modes and incorporate the effects of the small- and fast-scale
 1105 modes using extra damping and stochastic forcing in the approximate model. This follows
 1106 the basic idea of the stochastic mode reduction strategy^{53–56}, though the manipulation here
 1107 is less sophisticated. It is also important to notice that the extra stochastic noise added into
 1108 the approximate model is crucial since the energy in modes ψ_k for $k = 3, \dots, 10$ as shown
 1109 in Panels (e) and (p) of Figure 19 is non-negligible. Without these extra stochastic noise,
 1110 the total variance will be underestimated, which will severely affect the prediction skill of
 1111 the extreme events in the system.

1112 For the simplicity of notation, we make a change of variables,

$$\psi_1 = \frac{1}{2\sqrt{2}} \left((v_2 - v_1) - (v_2 + v_1)i \right), \quad \text{and} \quad \psi_2 = \frac{1}{2\sqrt{2}} \left((v_4 - v_3) - (v_4 + v_3)i \right). \quad (34)$$

1113 and therefore the 5-mode approximate model is given by,

$$\begin{aligned} \frac{du}{dt} &= \omega_1 v_1 + 2\omega_3 v_3 - d_u u + \sigma_u \dot{W}_u, \\ \frac{dv_1}{dt} &= -\beta v_2 + v_2 u - 2\omega_1 u - d_{v_1} v_1 + \sigma_1 \dot{W}_1, \\ \frac{dv_2}{dt} &= \beta v_1 - v_1 u - d_{v_2} v_2 + \sigma_2 \dot{W}_2, \\ \frac{dv_3}{dt} &= -\frac{\beta}{2} v_4 + 2v_4 u - \omega_3 u - d_{v_3} v_3 + \sigma_3 \dot{W}_3, \\ \frac{dv_4}{dt} &= \frac{\beta}{2} v_3 - 2v_3 u - d_{v_4} v_4 + \sigma_4 \dot{W}_4, \end{aligned} \quad (35)$$

1114 where $\omega_1 = H_1/\sqrt{2}$ and $\omega_3 = H_2/\sqrt{2}$. In (35), all the variables u, v_1, v_2, v_3 and v_4 are real.
 1115 The damping and stochastic forcing here are different from the perfect model since they now
 1116 also include some effects from the smaller scale modes of the perfect model that are ignored
 1117 here. Notably, the approximate model (35) satisfies the physics constraint, where the total
 1118 energy in the nonlinear terms is conserved^{35,36}.

1119 **C. Parameter estimation, data assimilation and long-term prediction skill**

1120 The system in (35) is a nonlinear system. In practice, the observational data of the leading
 1121 a few stream functions can be obtained. Therefore, we assume here the observational time
 1122 series of ψ_1 and ψ_2 are available. As in many real applications of atmosphere and ocean, the
 1123 observational training data is very limited. Here only the short period as shown in Panel
 1124 (a) or Panel (l) of Figure 19 is used for model calibration. On the other hand, we assume
 1125 that there is no observations for the zonal flow u . Recall that u plays an important role in
 1126 transferring energy with the small-scale stream functions in a nonlinear way and altering the
 1127 system between zonally blocked and unblocked patterns. Thus, for predicting the extreme
 1128 events in the system, assimilating the unobserved zonal flow u becomes necessary. Note
 1129 that despite the intrinsic nonlinearity in the coupled system (35), the system belongs to the
 1130 conditional Gaussian framework as was discussed in Section II, which allows an efficient way
 1131 of implementing parameter estimation and data assimilation.

1132 **Parameter estimation.**

1133 Applying the parameter estimation algorithm described in Section III D, we arrive at the
 1134 following estimated parameters in the approximate model (35),

$$\begin{aligned}
 \text{Regime } p = 1 : \quad & d_u = 0.0132, & d_v = 0.0187, & \sigma_u = 0.0515, & \sigma_v = 0.0501, \\
 & \omega_1 = 0.7035, & \omega_3 = 0.3508, & \beta = 1.9954, & \\
 \text{Regime } p = 0.5 : \quad & d_u = 0.1417, & d_v = 0.0205, & \sigma_u = 0.1450, & \sigma_v = 0.0504, \\
 & \omega_1 = 0.6712, & \omega_3 = 0.3485, & \beta = 1.9963, & \\
 & & & & (36)
 \end{aligned}$$

1135 where we have assumed all the damping coefficients in the v_i equations are the same and
 1136 all equal to d_v . Similar assumption is used for the stochastic forcing coefficients in the v_i
 1137 equations which all equal to σ_v .

1138 For the estimated parameters, those with clear physical meanings, for example β, ω_1 and
 1139 ω_3 , are quite close to the truth. The other parameters, mainly the stochastic forcing and
 1140 damping coefficients, are different from those in the perfect model. Note in particular that
 1141 the noise coefficients in the approximate model are larger than those in the perfect model.
 1142 Such a judicious model error with noise inflation compensates the error in the approximate
 1143 model due to the ignorance of the small-scale stream functions ψ_k from $k = \pm 3$ to ± 10 .

1144 **Data assimilation.**

1145 Using the approximate model (35) as the forecast model for data assimilation of the
1146 zonal flow u , the assimilated values are almost the same as the truth (figures not shown
1147 here) with a pattern correlation between the truth and the posterior mean states being 0.98
1148 and 0.95 in $p = 1$ and $p = 0.5$ regimes, respectively. In addition, the amplitudes of the
1149 assimilated states and the truth are also comparable with each other, implying the success
1150 of assimilating the extreme events. These results indicate the skill of using the approximate
1151 model in real-time state estimation of the unobserved process and accurately recovering the
1152 overall flow structure.

1153 **Long-range forecast.**

1154 Figure 20 shows the long-term forecast results. Panels (a)–(b) present model trajectories
1155 of ψ_1, ψ_2 and u simulated from the perfect model (28) and the approximate model (35)
1156 in $p = 1$ regime. These are simply a free run of each model and therefore we do not
1157 expect point-to-point correspondence between the two simulations due to the randomness.
1158 Nevertheless, these trajectories indicate that the qualitative features from both the models
1159 are similar. In particular, the approximate model succeeds in recovering the regime switching
1160 behavior in u . In Panels (c)–(d), the ACFs and PDFs associated with both the models are
1161 illustrated. The approximate model is quite skillful in capturing the strong oscillation, weak
1162 oscillation and the slowly but non-exponential decay in the ACFs associated with ψ_1, ψ_2 and
1163 u respectively. The approximate model also succeeds in recovering the PDFs of all the three
1164 variables, especially the variance which is important for predicting the extreme events in
1165 short- and medium-range, as will be discussed in the next subsection. Similar conclusions
1166 can be made in $p = 0.5$ regime. The only slight error lies in tracing the fast decay ACF of
1167 u in the approximate model. But the equilibrium PDFs and the ACFs associated with ψ_1
1168 and ψ_2 are recovered with high accuracy.

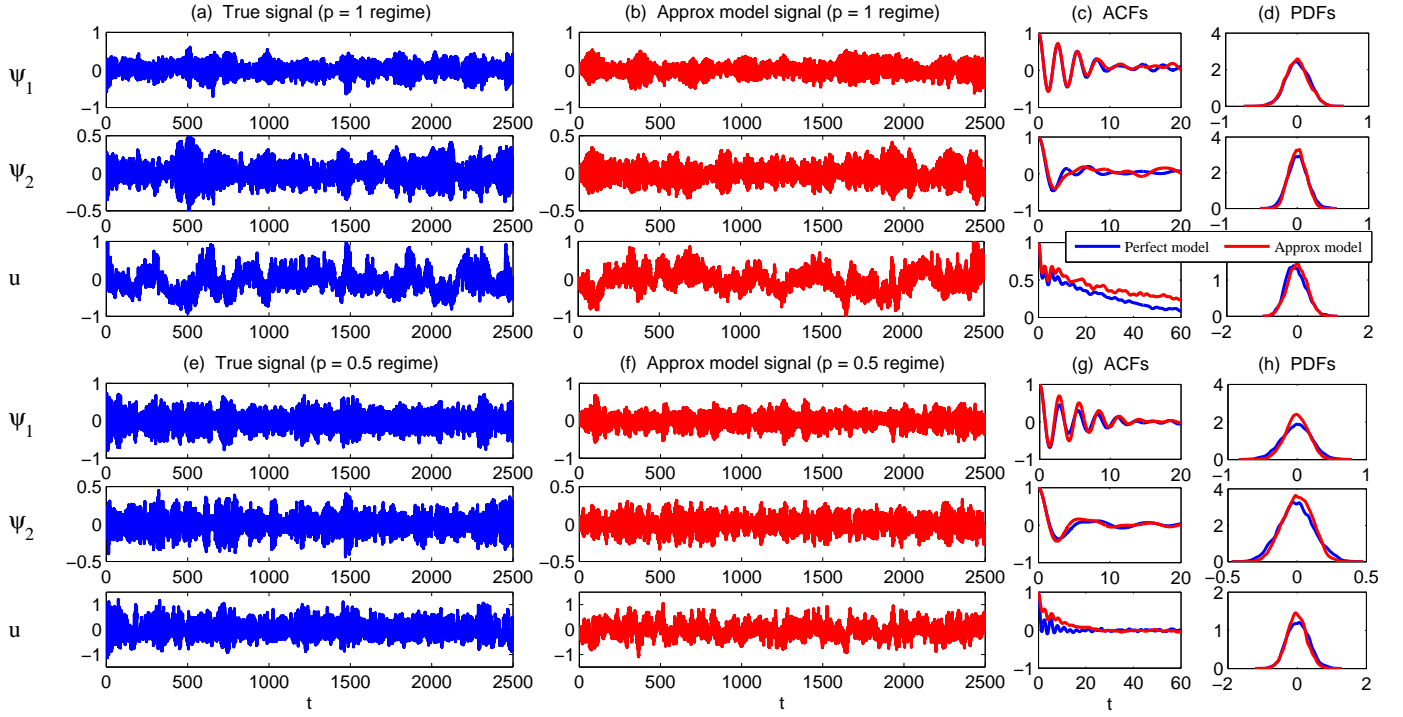


FIG. 20. Long-range forecasts of the layered topographic model. Panels (a)–(d): trajectories of the perfect model, trajectories of the approximate model, ACFs and PDFs in $p = 1$ regime. Panels (e)–(h): similar but for $p = 0.5$ regime.

1169 D. Short- and medium-range forecast

1170 With the approximate model and the assimilated initial conditions of u , the ensemble
1171 forecast is applied to study the short- and medium-range forecasts.

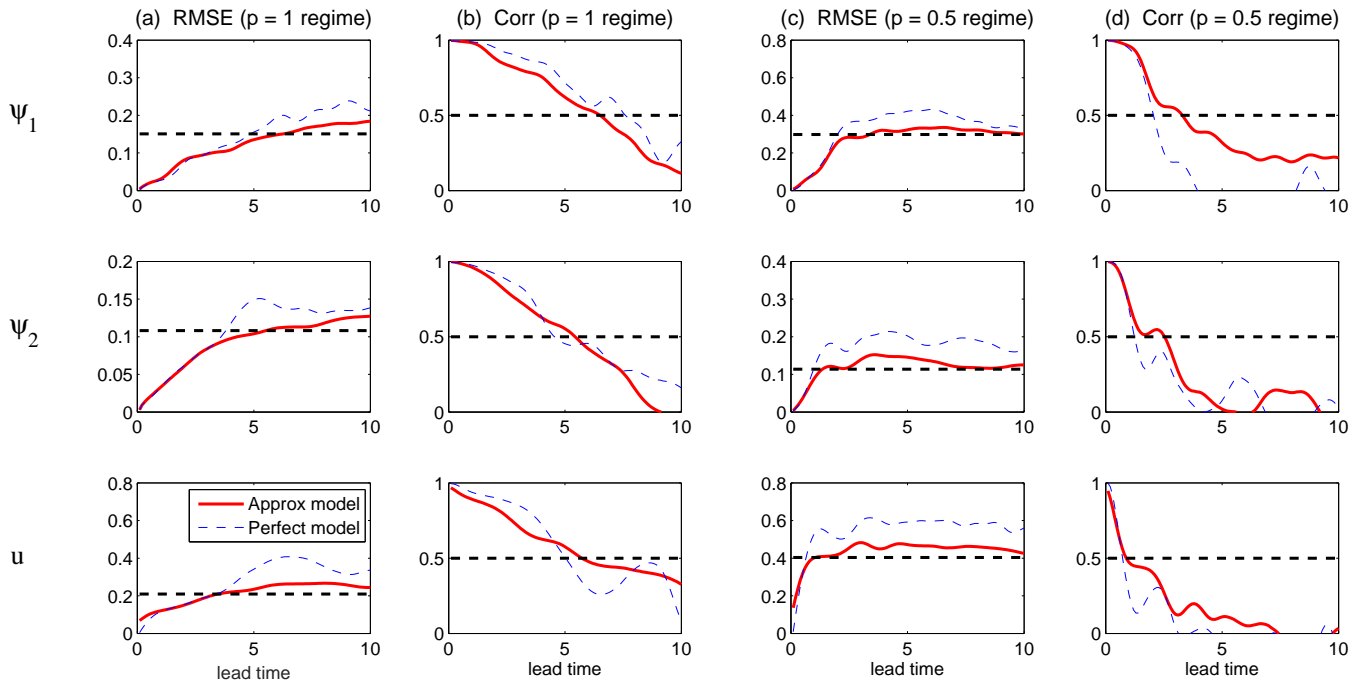
1172 Figure 21 shows the RMSE and Corr in the ensemble mean forecast as a function of the
1173 lead time. As comparison, the prediction using the perfect model is also included (blue). The
1174 approximate model has essentially the same skill as the perfect model in predicting all the
1175 three variables ψ_1, ψ_2 and u . The useful prediction based on these path-wise measurements
1176 as well as the information criterion for comparing the predicted amplitudes (not shown here)
1177 in $p = 1$ regime is about 5 units for all the three variables and that in $p = 0.5$ regime is
1178 3.5, 2.5 and 1 units for ψ_1, ψ_2 and u , respectively. Figure 22 shows the predicted trajectories
1179 at lead time 1, 2 and 3 units. The prediction of the extreme events up to 3 lead time units
1180 in $p = 1$ is quite accurate in terms of both the predicted patterns and the amplitudes. The
1181 $p = 0.5$ regime has a shorter range of useful predictions, but the overall skill up to 1 unit
1182 for both quiescent and extreme events are significant.

1183 Some case studies are included in Figure 23. In Panels (a)–(c), the ensemble prediction
1184 starts from $t = 300, 1390$ and 1460 , respectively, and each prediction is run for 30 units
1185 forward. Although the overall skillful prediction in $p = 1$ regime as shown in Figure 21 is
1186 5 units, the three events in Panels (a)–(c) of Figure 23 indicate that the useful prediction
1187 depends on the initial phase and the follow-up structure of the signal. Despite the intrinsic
1188 chaotic behavior, the useful prediction in case study 1 reaches 12 units, where all the extreme
1189 events within this time interval are captured accurately by the approximate model. On the
1190 other hand, the prediction in case study 2 is completely unskillful due to the fact that u has
1191 no internal oscillation structure for this particular event while the long-term trend cannot
1192 be captured by the ensemble mean forecast. Case study 3 shows a skillful prediction up
1193 to 6 units where again the extreme events within this time interval are captured with high
1194 accuracy.

1195 Panels (d)–(e) in Figure 23 compare the predicted stream functions using the approximate
1196 model with the truth in the case study 1 from Panel (a) at lead times 1.5 and 6.3, where the
1197 truth is generated from the perfect 21-mode model. The true values of the large-scale zonal
1198 flow at these two time instants are $u = 0.727$ and $u = -0.03$, respectively. The approximate
1199 model is quite skillful in predicting the overall flow patterns. In particular, the predictions

1200 succeed in capturing the regime switching phenomenon with a zonally unblocked structure
 1201 at $t = 301.5$ and a zonally blocked structure at $t = 306.3$. There are some small errors in the
 1202 prediction. For example, in Panel (d) the true signal at $x = 0.8$ has a sudden increase in the
 1203 meridional velocity while it is missed in the prediction using the 5-mode approximate model.
 1204 This meridional velocity is actually triggered by the modes $\psi_k, k = 3, 4, \dots$, which are not
 1205 included approximate model. Therefore, even if u, ψ_1 and ψ_2 are predicted almost perfectly,
 1206 which is the case here, there can be a small intrinsic barrier in recovering the original field
 1207 because of the simplification of the model by dropping the smaller scale modes.

1208 Similar conclusions are reached for $p = 0.5$ regime, as can be see in Panels (f)–(j) in Figure
 1209 23, despite that the useful prediction becomes shorter due to the more intrinsic turbulent
 1210 behavior in this regime.



1211

1212 FIG. 21. Short- and medium-range forecasts of the layered topographic model. Panels (a)–(b)
 1213 show the RMSE and Corr as a function of the lead time in $p = 1$ regime. Panels (c)–(d) show
 1214 those in $p = 0.5$ regime. The black dashed lines in the RMSE panels show the standard deviation
 1215 of the true signal and those in the Corr panels show the Corr = 0.5 threshold.

1217

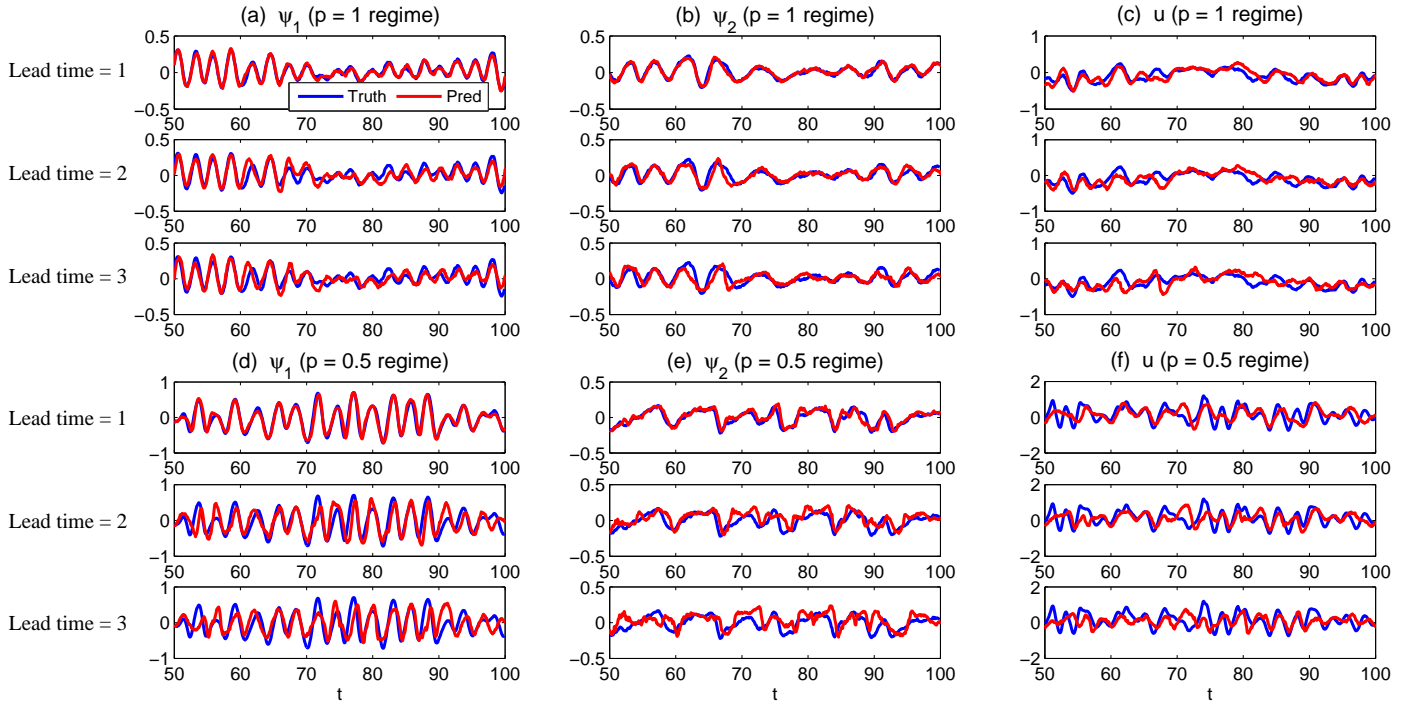


FIG. 22. The truth and the ensemble mean prediction at three different lead times 1, 2 and 3 using the approximate model with assimilated initial conditions of u . Panels (a)–(c): $p = 1$ regime. Panels (d)–(f): $p = 0.5$ regime.

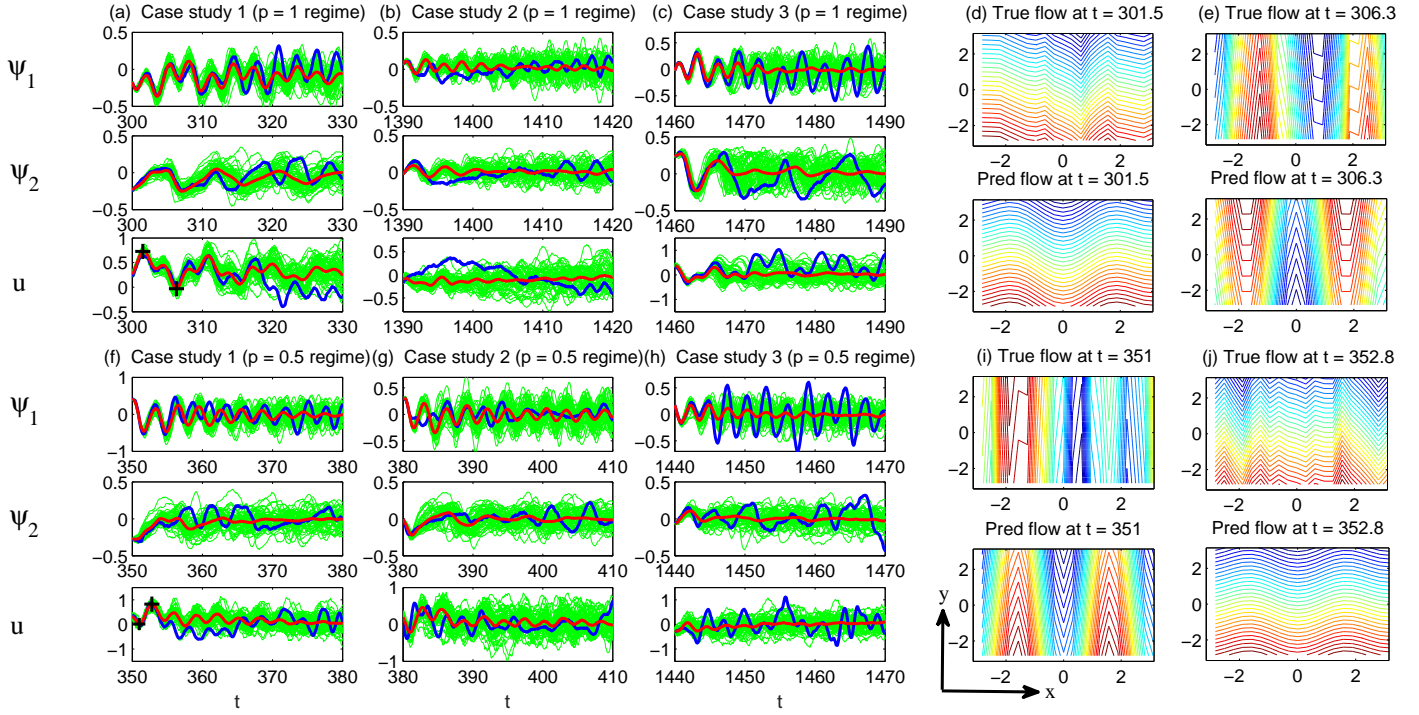
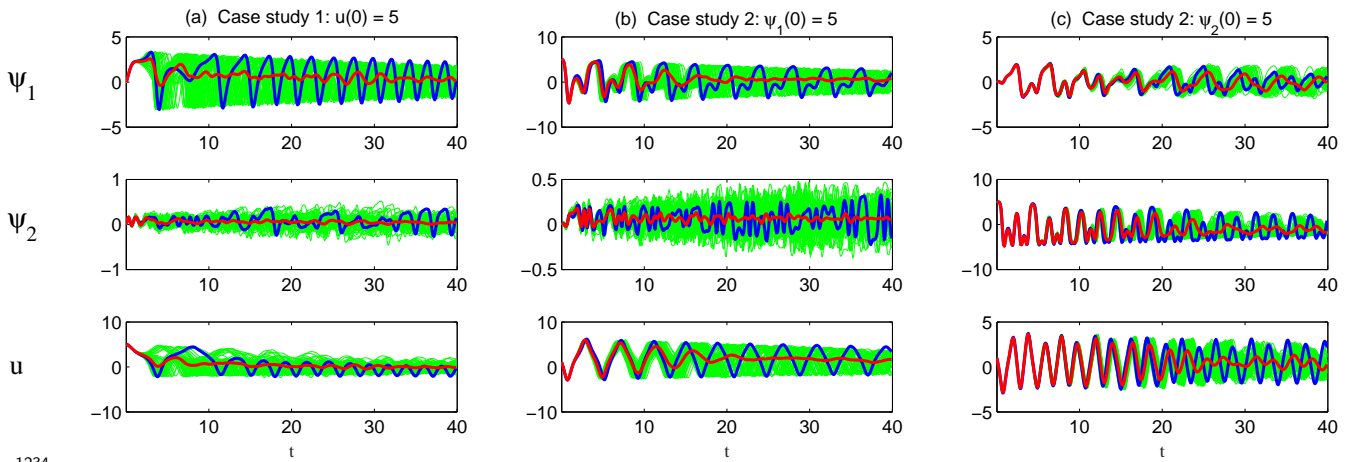


FIG. 23. Case studies of the ensemble prediction. Panels (a)–(e): $p = 1$ regime. Panels (f)–(j): $p = 0.5$ regime. In each subpanel of (a)–(c) and (f)–(h) the blue curve shows the truth and the red one shows the ensemble mean forecasts which is averaged over 50 ensemble members in green. Panels (d)–(e) compare the truth and the predicted overall streamlines in $p = 1$ regime at $t = 301.5$ and $t = 306.3$ (marked in black '+' in Panel (a)), where the starting time is $t = 300$. Panels (i)–(j) compare the truth and the predicted overall streamlines in $p = 0.5$ regime at $t = 351$ and $t = 352.8$ (marked in black '+' in Panel (f)), where the starting time is $t = 350$.

1218 **E. Prediction with an initial value starting outside the attractor**

1219 Finally, we study the prediction skill of the approximate model, which starts from a
 1220 value that is outside the attractor. In the three columns of Figure 24, we show the ensemble
 1221 predictions in $p = 1$ regime by letting the initial value of u , ψ_1 and ψ_2 be outside the attractor,
 1222 respectively. It is not difficult to tell that the true trajectories starting from values outside
 1223 the attractor behave in a very different way from those inside the attractor.

1224 When $u(0)$ is outside the attractor (Column (a)), the ensemble mean prediction using
 1225 the approximate model is accurate up to 3 units. The ensemble spread is very skillful in
 1226 capturing the envelope of the true signals as time evolves. When $\psi_1(0)$ is outside the attractor
 1227 (Column (b)). The useful ensemble mean prediction using the approximate model is about
 1228 10 units. The extreme events within this 10-unit interval in the zonal velocity are accurately
 1229 captured. Again, the ensemble spread clearly and accurately indicates the amplitudes in the
 1230 true signal. When $\psi_2(0)$ is outside the attractor (Column (c)). The skillful ensemble mean
 1231 prediction using the approximate model extends to 20 units! Note that within the first 10
 1232 units, the ensemble spread is very narrow, indicating the high confidence in the ensemble
 1233 mean prediction, including all the extreme events.



1234
 1235 FIG. 24. Prediction using the approximate model starting from a point that is outside the attractor.
 1236 Column (a) shows that when generating the true signal from the perfect model, the initial value of
 1237 ψ_1 is outside the attractor. Columns (b) and (c) show that the initial value of ψ_2 and u are outside
 1238 the attractor respectively. Again, the blue curve is the true and the red curve is the ensemble mean
 1239 prediction with the 50 ensemble members shown in green. Here $p = 1$.

1240 **VIII. A 6-DIMENSIONAL LOW-ORDER MODEL MIMICKING THE**
 1241 **CHARNEY-DEVORE (CDV) MODEL**

1242 **A. The perfect model and its properties**

1243 Charney and DeVore (CDV) made an fundamental contribution for the regime switching
 1244 behavior of the atmosphere¹²⁰. In this section, a 6-dimensional low-order model that mimics
 1245 the dynamical behavior of the CDV model is used as the perfect model. Despite the regime
 1246 switching behavior, this model has distinct mathematical structures and physical mecha-
 1247 nisms compared with the one studied in the previous section. It also possesses some unique
 1248 features, as will be discussed at the end of this subsection, that provide a very tough test
 1249 for predicting the extreme events and the transition behavior. The goal here is to design
 1250 suitable and efficient strategies of developing an approximate model that is able to predict
 1251 the extreme events and other non-Gaussian features in such a model.

1252 This 6-dimensional low-order model is obtained by a Galerkin projection and truncation
 1253 of the barotropic vorticity equation on a β -plane channel^{133,134}. The barotropic vorticity
 1254 equation is the following,

$$\frac{\partial}{\partial t} \nabla^2 \psi = -J(\psi, \nabla^2 \psi + f + \gamma h) - C \nabla^2 (\psi - \psi^*), \quad (37)$$

where the domain of longitude and latitude (x, y) are given by $[0, 2\pi] \times [0, \pi b]$. The parameter $b = 2B/L$ determines the ratio between the dimensional zonal length L and the meridional width B of the channel. The stream function ψ is periodic in x . The meridional boundaries $y = 0$ and $y = \pi$ have the conditions $\partial\psi/\partial x = 0$ and $\int_0^{2\pi} (\partial\psi/\partial y) dx = 0$. The Coriolis parameter f generates the beta effect in model. Orography enters with h , the orographic height, and is scaled with γ . J is the Jacobi operator and the damping coefficient C is the newtonian relaxation to the streamfunction profile ψ^* , which represents the forcing associated with the two zonal modes as will be discussed shortly. Next, the barotropic vorticity equation (37) is projected on a set of basis functions which are eigenfunctions of the Laplace operator ∇^2 ,

$$\phi_{0m}(y) = \sqrt{2} \cos(my/b), \quad \phi_{nm}(x, y) = \sqrt{2} e^{inx} \sin(my/b),$$

The 6-dimensional model is obtained by truncating the expansion of the stream function and the topographic height after $|n| = 1$ and $m = 2$. Then the time-dependent complex

variables of the stream functions $\psi_{01}, \psi_{02}, \psi_{\pm 11}, \psi_{\pm 12}$ are transformed to real variables:

$$\begin{aligned} x_1 &= \frac{1}{b}\psi_{01}, & x_2 &= \frac{1}{b\sqrt{2}}(\psi_{11} + \psi_{-11}), & x_3 &= \frac{i}{b\sqrt{2}}(\psi_{11} - \psi_{-11}), \\ x_4 &= \frac{1}{b}\psi_{02}, & x_5 &= \frac{1}{b\sqrt{2}}(\psi_{12} + \psi_{-12}), & x_6 &= \frac{i}{b\sqrt{2}}(\psi_{12} - \psi_{-12}), \end{aligned}$$

while the topography h is chosen to have only the (1, 1) wave profile,

$$h(x, y) = \cos(x) \sin(y/b).$$

1255 These manipulations lead to a 6-dimensional ODE model, where x_1, x_4 represent the zonal
1256 flow, x_2, x_3 are the topographic Rossby waves and x_5, x_6 are the Rossby waves.

1257 In the study here, extra small noise is added to this model, which allows some effects
1258 from the small-scale modes to enter into this low-order model. The noisy version of the
1259 6-dimensional CDV model reads,

$$\begin{aligned} dx_1 &= \left(\gamma_1^* x_3 - C(x_1 - x_1^*) \right) dt + \sigma_1 dW_1, \\ dx_4 &= \left(\gamma_2^* x_6 - C(x_4 - x_4^*) + \epsilon(x_2 x_6 - x_3 x_5) \right) dt + \sigma_4 dW_4, \\ dx_2 &= \left(-(\alpha_1 x_1 - \beta_1) x_3 - C x_2 - \delta_1 x_4 x_6 \right) dt + \sigma_2 dW_2, \\ dx_3 &= \left((\alpha_1 x_1 - \beta_1) x_2 - \gamma_1 x_1 - C x_3 + \delta_1 x_4 x_5 \right) dt + \sigma_3 dW_3, \\ dx_5 &= \left(-(\alpha_2 x_1 - \beta_2) x_6 - C x_5 - \delta_2 x_4 x_3 \right) dt + \sigma_5 dW_5, \\ dx_6 &= \left((\alpha_2 x_1 - \beta_2) x_5 - \gamma_2 x_4 - C x_6 + \delta_2 x_4 x_2 \right) dt + \sigma_6 dW_6. \end{aligned} \tag{38}$$

1260 Here the terms multiplied by α_i model the advection of the waves by the zonal flow. The
1261 β_i terms are due to the Coriolis force; the γ terms are generated by the topography. The
1262 C terms are the Newtonian damping to the zonal profile $x^* = (x_1^*, 0, 0, x_4^*, 0, 0)$. The δ -
1263 and ϵ -terms describe the nonlinear triad interaction between the zonal (0, 2) mode and the
1264 (1, 1) and (1, 2) waves. This triad is responsible for the possibility of barotropic instability
1265 of the (0, 2) mode. Note that the model is scaled such that 1 time unit in the model roughly
1266 corresponds to 1 day.

1267 Following^{133,134}, the following parameter values are taken: $C = 0.1$, corresponding to a
1268 damping time of 10 days; $\beta = 1.25$, corresponding to a channel centered around a latitude
1269 of 45° ; $b = 0.5$, the north-south extent of the channel is 25% of its east-west extent; and
1270 $x_1^* = 0.95$ and $x_4^* = -0.76095$. These parameters allow a combination of topographic and
1271 barotropic instabilities. The noise coefficients added here are $\sigma_1 = \dots = \sigma_6 = 0.005$. Such

1272 a choice of the noise coefficients allow the dynamical behavior of this stochastic model to
1273 remain similar to its deterministic version as was studied in^{133,134}.

1274 Note that x_1 and x_4 associated with ψ_{01} and ψ_{02} describe the zonal flows, and the forcing
1275 x_1^* and x_4^* are only imposed on these modes. Therefore, it is natural to assume x_1 and x_4
1276 are the observed variables while x_2, x_3, x_5 and x_6 are unobserved. The goal is to predict the
1277 extreme events in the system given only short trajectories of x_1 and x_4 .

1278 **Model properties.**

1279 Panels (a)–(c) of Figure 25 show the chaotic trajectories, non-Gaussian PDFs and the
1280 ACFs of the model (38). It is easy to tell from the model trajectories that this model
1281 has multiple equilibria, which is confirmed by the phase plot (x_1, x_4) in Panel (d) (two
1282 stable equilibria; top left and bottom right). The spatial patterns associated with these two
1283 equilibria are quite different with each other, as shown in Panels (e) and (f) for two sample
1284 events corresponding to the time instants marked in red in Panel (a) that lie near these two
1285 equilibria. The streamlines shown in Panel (e) corresponds to an equilibria with largely zonal
1286 character with strong zonal jets while that in Panel (f) is dominated by topographically
1287 effects with vortices and meander jets. When the blocking events happen, x_1 reaches its
1288 maximum while x_4 lies in its minimum value.

1289 Panel (b) shows the equilibrium PDFs of all the 6 model variables. The profiles of these
1290 PDFs are quite different: x_1 and x_3 have weakly bimodal distributions; x_2, x_4 and x_6 are
1291 highly skewed with an one-sided fat tail towards the negative side; and x_5 is skewed with
1292 a fat tail towards the positive side. Panel (c) illustrates the ACFs, which imply multiple
1293 decorrelation time scales of the system, where x_1, x_3 and x_4 has a much longer memory than
1294 x_2, x_5 and x_6 .

1295 One very interesting and important feature of this model (or more precisely its determinis-
1296 tic version) is that projecting this 6-dimensional model to its leading 5 Empirical Orthogonal
1297 Functions (EOFs) explains 99.5% of the variance. However, such a 5-dimensional projected
1298 dynamics completely misses the dynamical features in the original model, where the multiple
1299 equilibria disappears and the 5-dimensional model cannot reproduce regime transitions¹³⁴.
1300 Therefore, this 6-dimensional model provides a very useful and tough testbed for developing
1301 suitable approximate models in predicting the transition behavior and extreme events in
1302 highly chaotic systems.

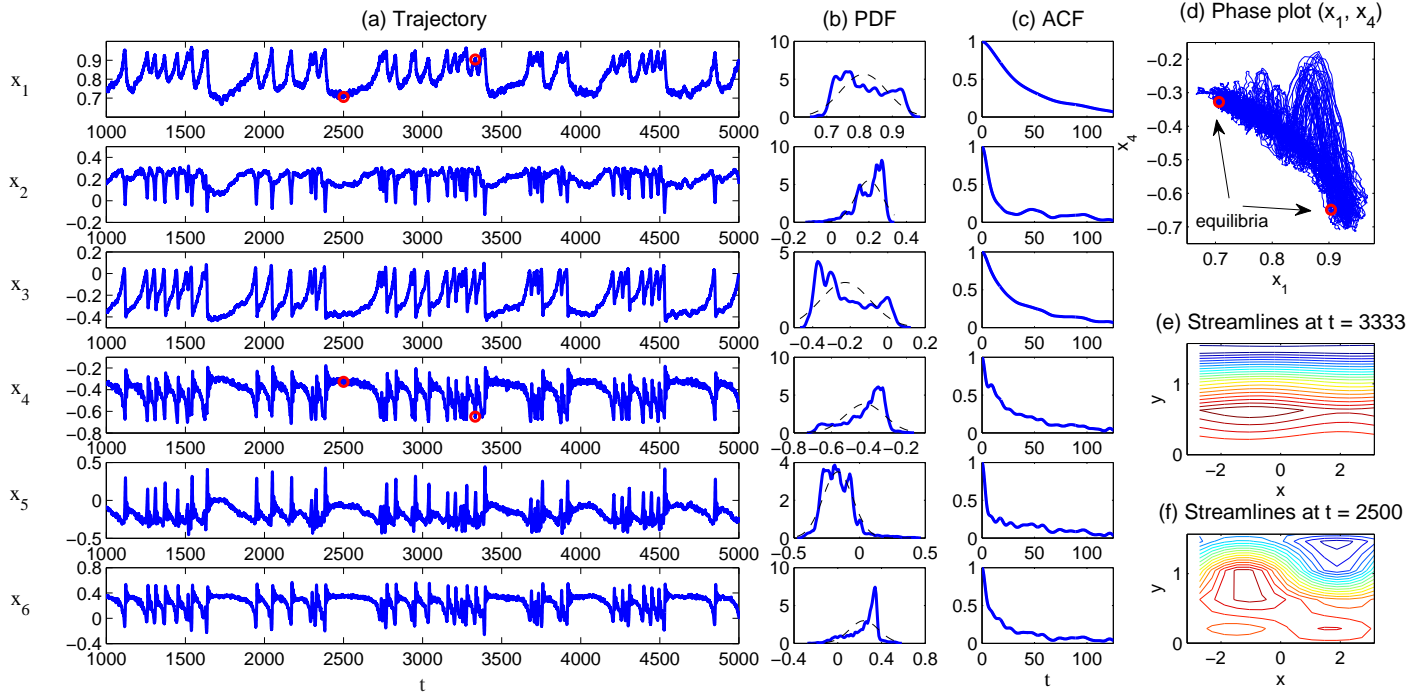


FIG. 25. Panels (a)–(c): Model trajectories, PDFs and ACFs of the 6-D CDV model (38). The black dashed lines in column (b) show the Gaussian fits of the PDFs. Panel (d): Phase plot of (x_1, x_4) . Panels (e)–(f): Streamlines at $t = 3333$ and $t = 2500$, corresponding to the time instants marked in red dots in Panel (a).

1303 **B. The approximate model**

1304 Our goal here is to develop a suitable approximate model for describing and predicting
1305 the key features of the 6-dimensional low-order CDV model (38). Recall that the conditional
1306 Gaussian nonlinear models in Section II allow an efficient and accurate data assimilation al-
1307 gorithm, which facilitates effective predictions. Therefore, it is natural to develop a suitable
1308 approximate model that belongs to the conditional Gaussian nonlinear modeling framework.
1309 Note that by observing x_1 and x_4 , the 6-dimensional CDV model (38) is not a conditional
1310 Gaussian model due to the nonlinear coupling term $\epsilon(x_2x_6 - x_3x_5)$. In fact, the topographic
1311 Rossby waves x_2, x_3 and the Rossby waves x_5, x_6 interact with each other through the above
1312 nonlinear coupled term. Only in the absence of the Rossby waves x_5, x_6 , the coupled sys-
1313 tem x_1, \dots, x_4 is a conditional Gaussian system. Therefore, suitable strategies need to be
1314 developed to cope with this nonlinear term in the approximate model.

1315 **Strategy 1: A bare truncation model.**

1316 The simplest way to deal with this nonlinear term is to build a bare truncation model,
1317 where the nonlinear term $\epsilon(x_2x_6 - x_3x_5)$ is completely dropped. However, this bare trunca-
1318 tion model suffers from finite time blowup issue. In fact, the blowup occurs very quickly and
1319 even for a very short lead time (much shorter than the decorrelation time), the predicted
1320 values have a large chance to go to infinity.

1321 **Strategy 2: A nonlinear approximate model with linear feedback terms.**

1322 Another straightforward idea is to replace the quadratic term $\epsilon(x_2x_6 - x_3x_5)$ by a combi-
1323 nation of four linear terms $c_1x_2 + c_2x_6 - c_3x_3 - c_4x_5$. This approximation is better than the
1324 bare truncation model in the sense that the system will not blow up in a very short term.
1325 However, the predicted trajectories from this model still have a high probability to blow up
1326 in a finite time. In addition, the skillful prediction only lasts for very short time even if the
1327 predicted amplitude remains finite within that range.

1328 **Strategy 3: An approximate model with a stochastic forcing term.**

1329 Instead of using a deterministic and linear way to parameterize the nonlinear quadratic
1330 term, a new strategy is developed here, which involves using a simple stochastic forcing
1331 process b_1 to describe the effect of the quadratic term $\epsilon(x_2x_6 - x_3x_5)$. The approximate

1332 model reads,

$$\begin{aligned}
dx_1 &= \left(\gamma_1^* x_3 - C(x_1 - x_1^*) \right) dt + \sigma_1 dW_1, \\
dx_4 &= \left(\gamma_2^* x_6 - C(x_4 - x_4^*) + b_1 \right) dt + \sigma_4 dW_4, \\
dx_2 &= \left(-(\alpha_1 x_1 - \beta_1) x_3 - C x_2 - \delta_1 x_4 x_6 \right) dt + \sigma_2 dW_2, \\
dx_3 &= \left((\alpha_1 x_1 - \beta_1) x_2 - \gamma_1 x_1 - C x_3 + \delta_1 x_4 x_5 \right) dt + \sigma_3 dW_3, \\
dx_5 &= \left(-(\alpha_2 x_1 - \beta_2) x_6 - C x_5 - \delta_2 x_4 x_3 \right) dt + \sigma_5 dW_5, \\
dx_6 &= \left((\alpha_2 x_1 - \beta_2) x_5 - \gamma_2 x_4 - C x_6 + \delta_2 x_4 x_2 \right) dt + \sigma_6 dW_6, \\
db_1 &= \left(-d_b b_1 + \sigma_b b_2 + f_b \right) dt + \sigma_b dW_b, \\
db_2 &= \left(-d_b b_2 - \sigma_b b_1 + f_b \right) dt + \sigma_b dW_b.
\end{aligned} \tag{39}$$

1333 This is motivated by the SPEKF model^{48,49}, where a stochastic forcing is able to automati-
1334 cally learn the missing information on the fly via online data assimilation. Here, we adopt the
1335 simplest possible choice — a stochastic forcing b_1 driven by a Gaussian process. Note that
1336 two new processes b_1 and b_2 are actually incorporated into the approximate model. They to-
1337 gether form a linear stochastic oscillator while only b_1 gives feedback to the x_4 process. The
1338 reason to impose an “oscillated” forcing is that all the variables x_i have chaotic oscillator
1339 structures and so does the nonlinear term $\epsilon(x_2 x_6 - x_3 x_5)$. As will be seen below, with this
1340 cheap stochastic strategy, the approximate model is able to avoid finite time blowup issue
1341 and it provides surprisingly skillful predictions in both short and medium ranges. Notably,
1342 treating b_1 and b_2 as the extra unobserved variables, the resulting 8-dimensional nonlin-
1343 ear system in (39) is a conditional Gaussian nonlinear model, where only x_1 and x_2 have
1344 observations. The estimated parameters are given by:

$$F_b = 0.0081, \quad \omega_b = 0.6815, \quad d_b = 0.1339, \quad \sigma_b = 0.01326. \tag{40}$$

1345 Below the focus will be on the data assimilation and prediction skill using the approximate
1346 model in (39).

1347 C. Data assimilation

1348 Since the approximate model with the stochastic forcing (39) is a conditional Gaussian
1349 system, the data assimilation algorithm (3) provides an efficient state estimation of both

1350 the unobserved variables x_2, x_3, x_5, x_6 and the stochastic forcing b_1, b_2 , which are shown
1351 in Figure 26 (in red color). As comparison, we also show the truth of the unobserved
1352 variables x_2, x_3, x_5, x_6 (in blue color). It is clear that the data assimilation with the help
1353 of such a stochastic forcing term provides very accurate estimation of the hidden variables
1354 x_2, x_3, x_5, x_6 , where the pattern correlation of the assimilated and the true signals is higher
1355 than 0.95 for all the variables.

1356 Another striking result is presented in Panels (c) and (d) of Figure 26, where a comparison
1357 between the assimilated state of the stochastic forcing b_1 and the nonlinear term $\epsilon(x_2x_6 -$
1358 $x_3x_5)$ computed from the perfect model is illustrated. It is clear that the stochastic forcing
1359 b_1 almost perfectly recovers the nonlinear feedback, especially at the time instants that the
1360 nonlinear feedback is intermittent. This is a very important feature because it guarantees
1361 that the stochastic forcing is able to, at least for a short term, play the role of the nonlinear
1362 feedback term in prediction.

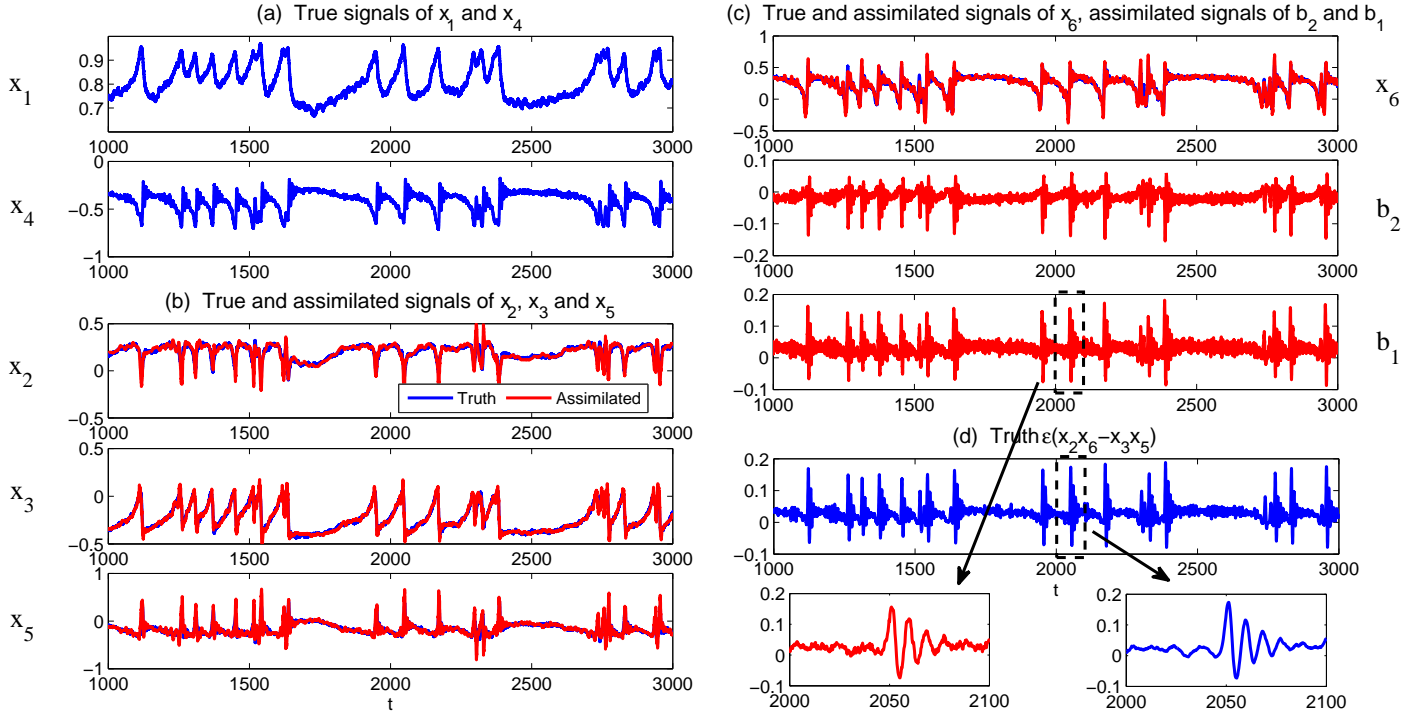


FIG. 26. Data assimilation of the 6-D CDV model using the approximate model (39). Panel (a): the true signals of x_1 and x_4 . Panels (b)–(c) the true and the assimilated posterior mean of x_2, x_3, x_5 and x_6 , and the assimilated stochastic forcing b_2 and b_1 . Panel (d): the true value of the nonlinear term $\epsilon(x_2x_6 - x_3x_5)$. A zoomed-in period of b_1 and the nonlinear term is also shown for comparison.

1363 **D. Predictions**

1364 **Short- and medium-range forecasts.**

1365 Our focus now is on the short- and medium-range forecasts. In Figure 27, the prediction
1366 skill in terms of the RMSE and Corr as a function of lead time is presented. The blue curves
1367 show the predictions using the perfect model with the perfect initial conditions; the red
1368 curves show those using the approximate model (39) with the perfect initial conditions; and
1369 the green curves show those using the approximate model (39) and the assimilated initial
1370 conditions. The ensemble mean is used here for computing the RMSE and Corr.

1371 Despite that the prediction using the approximate model (39) is less skillful than that
1372 using the perfect model as the increase of lead time, it is clear that the useful prediction for all
1373 the variables using the approximate model (39) is still at least 8 units. For some variables
1374 such as x_3 the useful prediction is 16 days and for x_1 it is much longer. In addition,
1375 the approximate model (39) using the assimilated initial conditions has nearly the same
1376 prediction skill as that using the perfect initial conditions, which verifies the accuracy in the
1377 assimilated states. These results imply that the approximate model is a suitable model for
1378 both short- and medium-range forecasts of such an extremely tough test model.

1379 Figures 28–29 include two case studies of the prediction tests. The ensemble mean pre-
1380 diction shown in Figure 28 is extremely accurate for both short and medium ranges. On the
1381 other hand, although the ensemble mean prediction in Figure 29 has a slight phase shift,
1382 which results in the deterioration of the pattern correlation, the overall prediction using the
1383 approximate model remains skillful. From these figures, it is clear that most of the extreme
1384 events take around 8 units to develop from the onset phase to the peak, which is within the
1385 skillful prediction range of the approximate model (39). Therefore, the approximate model
1386 is able to predict the entire development phase of the extreme events. On the other hand,
1387 starting from the peak of an extreme event, the approximate model succeeds in predicting
1388 the returning path to the quiescent state. In addition to the ensemble mean, the ensemble
1389 envelope also plays an important role in the prediction here. In fact, despite a slight phase
1390 shift in the ensemble mean prediction in Figure 29, the ensemble envelope clearly predicts
1391 the correct overall time evolution trends of the truth. Admittedly, the ensemble spread
1392 using the approximate model with the assimilated initial condition is larger than the pre-
1393 diction with the perfect initial condition, which is mainly due to the initial uncertainty in

1394 assimilating the hidden variables and the uncertainty introduced from the stochastic forcing.
1395 Nevertheless, the correct trends are still unambiguously predicted by the ensemble members
1396 in both short and medium terms.

1397 **Long range forecast.**

1398 The approximate model (39) fails to reproduce the same long-term equilibrium PDFs as
1399 the truth. This is not surprising since the stochastic forcing for the long range forecast loses
1400 its memory of the initial condition and essentially becomes a constant. Its contribution to
1401 the system is then quite different from the original nonlinear term $\epsilon(x_2x_5 - x_3x_6)$, which
1402 evolves in time. Note that the study in the previous work⁷⁴ has already illustrated that in
1403 the presence of model error it is extremely difficult to develop suitable approximate models
1404 that are able to simultaneously have both short and long range forecast skill. Nevertheless,
1405 the approximate model (39) is still able to provide some useful information for the long-
1406 range forecasts. First, the approximate model (39) avoids finite-time blowup issue and
1407 its equilibrium PDFs contain non-Gaussian features, which already outweighs many other
1408 approximation strategies, such as bare truncation and linear approximations, for describing
1409 strongly chaotic systems. Second, the ACFs of x_1, x_2 and x_3 from the approximate model
1410 (39) are quite similar to the truth and the errors in the ACFs of x_4, x_5 and x_6 are also only
1411 moderate. These features in the ACFs together with the accurate data assimilation results
1412 actually guarantee the skillful short- and medium-range forecasts.

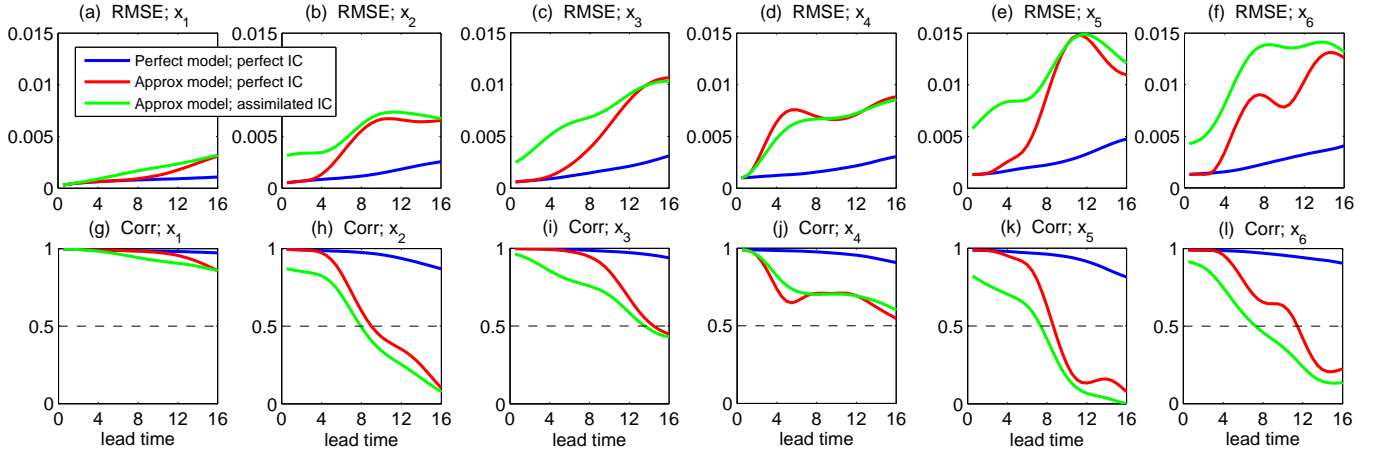


FIG. 27. Short- and medium-range forecasts using the perfect model and the approximate model (39). Top and bottom rows show the RMSE and Corr and a function of lead time. The blue curves show the predictions using the perfect model with the perfect initial conditions; the red curves show those using the approximate model (39) model with the perfect initial conditions; and the green curves show those using the approximate model (39) and the assimilated initial conditions. The ensemble mean is used here for computing the RMSE and Corr with the truth.

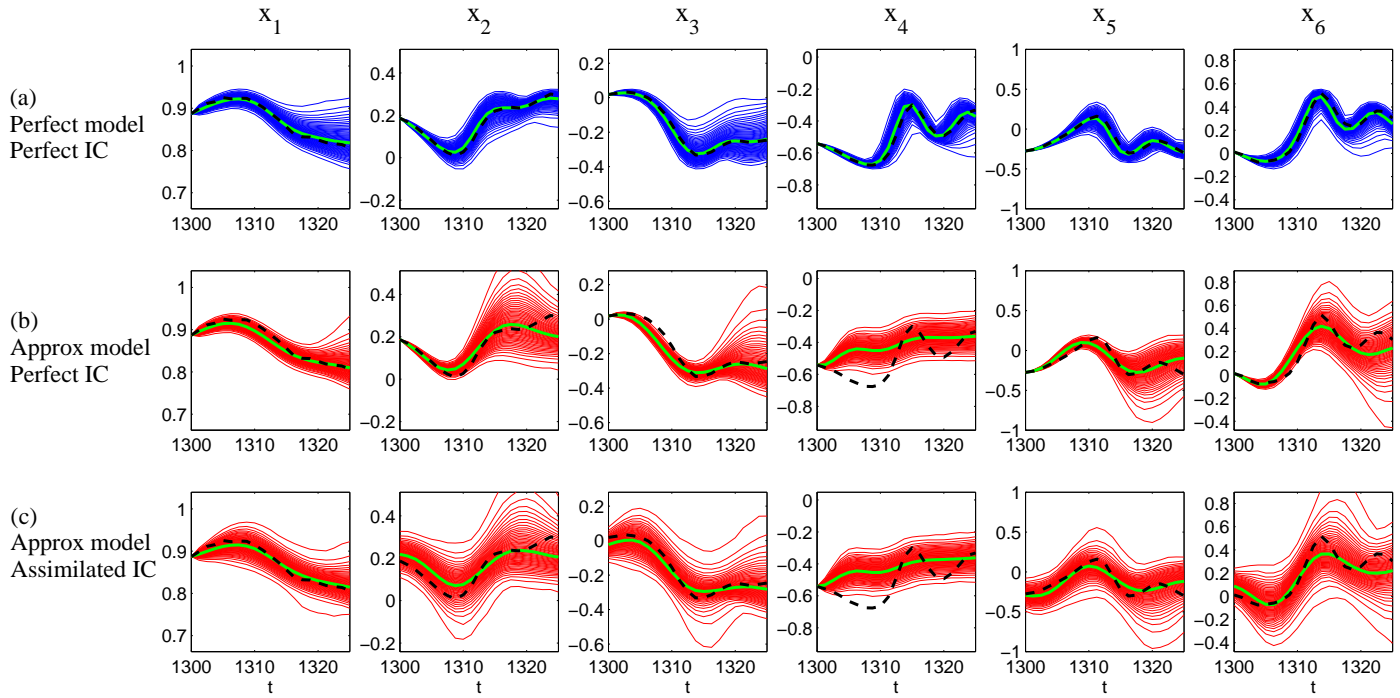


FIG. 28. Case study. Prediction starting from $t = 1300$. Note that each PDF is shown with 50 thin curves (blue for the perfect model and red for the approximate model), which represent the 1st, 3rd, 5th, ..., 97th and 99th percentiles of the of the PDF. The green curve represents the mode of the PDF since the PDF is non-Gaussian. The black dashed curve is the true signal.

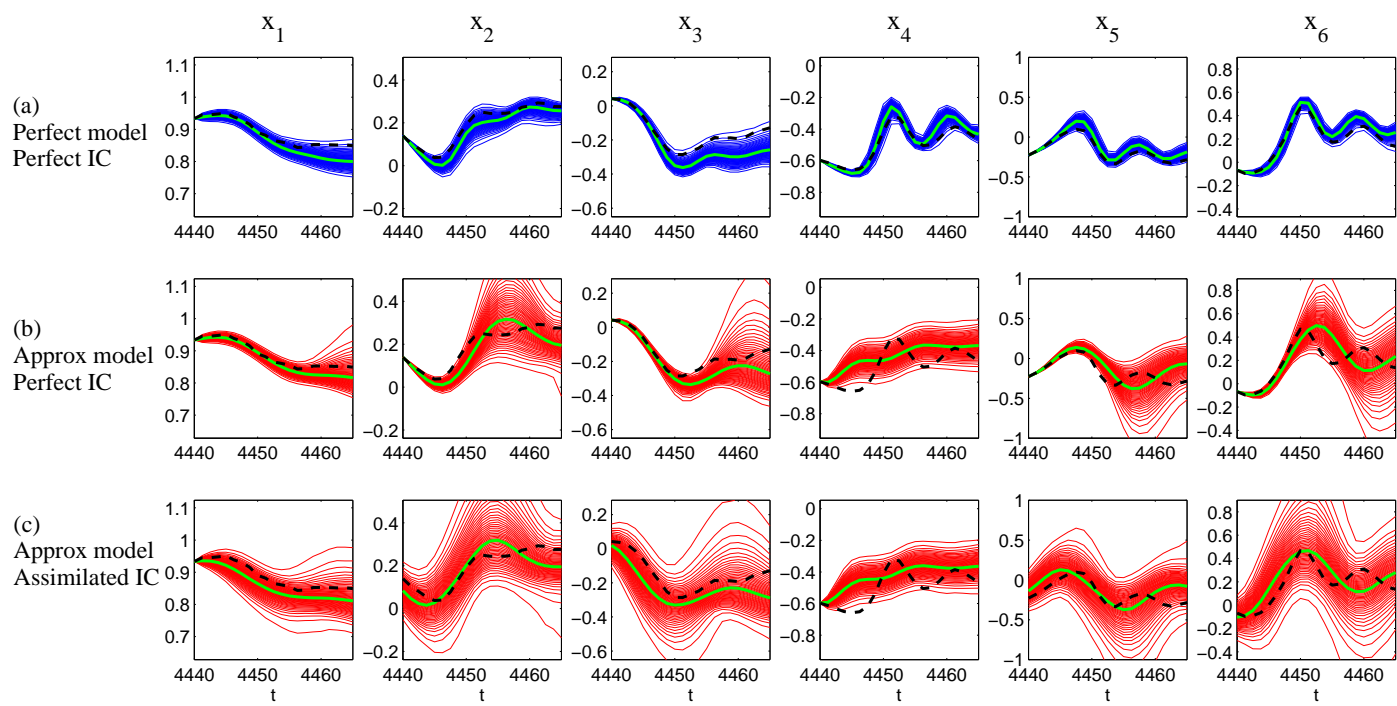


FIG. 29. Case study. Prediction starting from $t = 4440$.

1413 IX. CONCLUSION

1414 Extreme events appear in many complex nonlinear dynamical systems. Predicting ex-
1415 tremes events has both scientific significance and practical implications. The main difficulties
1416 in predicting the extreme events include the lack of a complete understanding of physics,
1417 the unaffordable computational cost of running the complex dynamical systems and the
1418 errors in data assimilation or state estimation. Notably, in many practical situations, only
1419 partially observed time series are available for model calibration and the training period is
1420 often very short. All these facts result in great challenges and lead to the failure of many
1421 purely data-driven methods in the extreme events prediction.

1422 In this paper, a new mathematical framework of building suitable nonlinear approximate
1423 models is developed, which aims at predicting both the observed and hidden extreme events
1424 in complex nonlinear dynamical systems using only short and partially observed training
1425 time series. The models belonging to this mathematical framework are highly nonlinear and
1426 are able to capture many key non-Gaussian characteristics as observed in nature. Physically
1427 motivated processes and physics constraints can be incorporated into the models, which make
1428 this framework fundamentally different from many purely data-driven statistical models that
1429 have no clear physical meanings. Such a feature also allows using only a short training time
1430 series for model calibration. In addition, this modeling framework provides closed analytic
1431 formulae for assimilating the states of the unobserved variables, which is computationally
1432 efficient and accurate. The details of this modeling framework is shown in Section II. Section
1433 III contains the efficient and accurate data assimilation, parameter estimation and prediction
1434 algorithms as well as the details of using both the path-wise and information measurements
1435 in quantifying the prediction skill. Different effective and practical strategies of developing
1436 suitable approximate models for predicting extreme events and other non-Gaussian features
1437 in various complex turbulent dynamical systems are illustrated in Section IV to Section
1438 VIII.

1439 In Section IV, the skill of applying a cheap stochastic parameterization to approximate the
1440 complicated dynamical behavior in the hidden process is explored. This simple and efficient
1441 stochastic parameterization is able to recover the nonlinear feedback from the unresolved
1442 variable to the observed one. Notably, the nonlinear approximate model with such a cheap
1443 stochastic parameterization has nearly the same skill in predicting the extreme events at

1444 all short, medium and long ranges. Section V makes use of a nonlinear dyad model to
1445 show the success of applying a simple feedback control strategy in the approximate model
1446 to facilitate the prediction of the hidden extreme events, which is a great challenge given
1447 only partial observations. In Section VI, the Lorenz 63 model is used as a simple test model
1448 for predicting extreme events in the intrinsic chaotic models. The goal for testing this
1449 model is to understand the model error due to the noise inflation in affecting the extreme
1450 events prediction, where the noise inflation is a typical strategy of developing approximate
1451 models in many real applications. It is shown that a moderate noise inflation retains the
1452 skill of the extreme events prediction at all short, medium and long ranges. Next, regime
1453 switching between multiple metastable states is a key feature in many nonlinear turbulent
1454 dynamical systems. Section VII starts with a 21-dimensional nonlinear topographic mean
1455 flow interaction model with regime switching. A simplified version of the stochastic mode
1456 reduction strategy is applied in a suitable way to develop an approximate physics-constrained
1457 nonlinear model with only 5 dimensions. The 5-dimensional physics-constrained nonlinear
1458 model has a significant skill in predicting both the observed and hidden extreme events as
1459 well as other non-Gaussian features, nearly the same as the perfect model prediction. It
1460 also succeeds in predicting the regime switching between the zonally blocked and unblocked
1461 patterns with high accuracy. In Section VIII, a 6-dimensional low-order Charney-DeVore
1462 (CDV) model is used as a testbed for predicting extreme events. This model is highly
1463 nonlinear and has strong chaotic features. The leading 5 EOFs contain more than 99.5%
1464 of the explained variance but they completely miss the nonlinear dynamical features and
1465 the regime switching behavior. Therefore, this 6-dimensional model is an extremely tough
1466 test model for predicting the intrinsic nonlinear transitions and extreme events. It is shown
1467 that a simple but judicious linear stochastic process with additive noise and memory has
1468 a significant skill in learning certain complicated nonlinear effects of this model on the fly.
1469 The resulting approximate nonlinear model by incorporating such a simple stochastic process
1470 allows efficient and accurate data assimilation. It succeeds in predicting both the observed
1471 and hidden extreme events in short and medium terms.

1472 **ACKNOWLEDGMENTS**

1473 The research of N.C. is supported by the Office of Vice Chancellor for Research and
1474 Graduate Education (VCRGE) at University of Wisconsin-Madison. The research of A.J.M.
1475 is partially supported by the Office of Naval Research Grant ONR MURI N00014-16-1-2161
1476 and the Center for Prototype Climate Modeling (CPCM) at New York University Abu
1477 Dhabi Research Institute. The research of both N.C. and A.J.M is supported by ONR
1478 MURI N00014-19-1-2421.

1479 **REFERENCES**

- 1480 ¹A. J. Majda, *Introduction to turbulent dynamical systems in complex systems* (Springer,
1481 2016).
- 1482 ²A. Majda and X. Wang, *Nonlinear dynamics and statistical theories for basic geophysical*
1483 *flows* (Cambridge University Press, 2006).
- 1484 ³S. H. Strogatz, *Nonlinear Dynamics and Chaos: with Applications to Physics, Biology,*
1485 *Chemistry, and Engineering* (CRC Press, 2018).
- 1486 ⁴D. Baleanu, J. A. T. Machado, and A. C. Luo, *Fractional dynamics and control* (Springer
1487 Science & Business Media, 2011).
- 1488 ⁵T. Deisboeck and J. Y. Kresh, *Complex systems science in biomedicine* (Springer Science
1489 & Business Media, 2007).
- 1490 ⁶J. Stelling, A. Kremling, M. Ginkel, K. Bettenbrock, and E. Gilles, *Foundations of*
1491 *Systems Biology* (MIT press, 2001).
- 1492 ⁷S. A. Sheard and A. Mostashari, “Principles of complex systems for systems engineering,”
1493 *Systems Engineering* **12**, 295–311 (2009).
- 1494 ⁸D. C. Wilcox, “Multiscale model for turbulent flows,” *AIAA journal* **26**, 1311–1320 (1988).
- 1495 ⁹E. Pelinovsky, C. Kharif, *et al.*, *Extreme ocean waves* (Springer, 2008).
- 1496 ¹⁰W. K.-M. Lau and D. E. Waliser, *Intraseasonal variability in the atmosphere-ocean climate*
1497 *system* (Springer Science & Business Media, 2011).
- 1498 ¹¹A. J. Clarke, *An introduction to the dynamics of El Niño and the Southern Oscillation*
1499 (Elsevier, 2008).
- 1500 ¹²A. J. Majda, M. Moore, and D. Qi, “Statistical dynamical model to predict extreme events

- 1501 and anomalous features in shallow water waves with abrupt depth change,” Proceedings
1502 of the National Academy of Sciences **116**, 3982–3987 (2019).
- 1503 ¹³J. D. Neelin, B. R. Lintner, B. Tian, Q. Li, L. Zhang, P. K. Patra, M. T. Chahine, and
1504 S. N. Stechmann, “Long tails in deep columns of natural and anthropogenic tropospheric
1505 tracers,” *Geophysical Research Letters* **37** (2010).
- 1506 ¹⁴A. J. Majda and B. Gershgorin, “Elementary models for turbulent diffusion with complex
1507 physical features: eddy diffusivity, spectrum and intermittency,” *Philosophical Trans-*
1508 *actions of the Royal Society A: Mathematical, Physical and Engineering Sciences* **371**,
1509 20120184 (2013).
- 1510 ¹⁵A. J. Majda and X. T. Tong, “Intermittency in turbulent diffusion models with a mean
1511 gradient,” *Nonlinearity* **28**, 4171 (2015).
- 1512 ¹⁶E. G. Altmann and H. Kantz, “Recurrence time analysis, long-term correlations, and
1513 extreme events,” *Physical Review E* **71**, 056106 (2005).
- 1514 ¹⁷W. L. Oberkampf, S. M. DeLand, B. M. Rutherford, K. V. Diegert, and K. F. Alvin,
1515 “Error and uncertainty in modeling and simulation,” *Reliability Engineering & System*
1516 *Safety* **75**, 333–357 (2002).
- 1517 ¹⁸M. Collins, B. B. Booth, B. Bhaskaran, G. R. Harris, J. M. Murphy, D. M. Sexton, and
1518 M. J. Webb, “Climate model errors, feedbacks and forcings: a comparison of perturbed
1519 physics and multi-model ensembles,” *Climate Dynamics* **36**, 1737–1766 (2011).
- 1520 ¹⁹T. Palmer and J. Räisänen, “Quantifying the risk of extreme seasonal precipitation events
1521 in a changing climate,” *Nature* **415**, 512 (2002).
- 1522 ²⁰A. J. Majda, “Challenges in climate science and contemporary applied mathematics,”
1523 *Communications on Pure and Applied Mathematics* **65**, 920–948 (2012).
- 1524 ²¹E. Kalnay, *Atmospheric Modeling, Data Assimilation and Predictability* (Cambridge uni-
1525 versity press, 2003).
- 1526 ²²A. Majda and N. Chen, “Model error, information barriers, state estimation and predic-
1527 tion in complex multiscale systems,” *Entropy* **20**, 644 (2018).
- 1528 ²³I. M. Navon, “Data assimilation for numerical weather prediction: a review,” in *Data*
1529 *assimilation for atmospheric, oceanic and hydrologic applications* (Springer, 2009) pp.
1530 21–65.
- 1531 ²⁴X. Zou, I. Navon, and F. Le Dimet, “Incomplete observations and control of gravity waves
1532 in variational data assimilation,” *Tellus A: Dynamic Meteorology and Oceanography* **44**,

- 1533 273–296 (1992).
- 1534 ²⁵W. Lahoz, B. Khatatov, and R. Ménard, “Data assimilation and information,” in *Data*
1535 *Assimilation* (Springer, 2010) pp. 3–12.
- 1536 ²⁶A. J. Majda and J. Harlim, *Filtering Complex Turbulent Systems* (Cambridge University
1537 Press, 2012).
- 1538 ²⁷G. Evensen, *Data Assimilation: the Ensemble Kalman Filter* (Springer Science & Business
1539 Media, 2009).
- 1540 ²⁸K. Law, A. Stuart, and K. Zygalakis, *Data Assimilation: a Mathematical Introduction*,
1541 Vol. 62 (Springer, 2015).
- 1542 ²⁹D. Qi and A. J. Majda, “Predicting extreme events for passive scalar turbulence in two-
1543 layer baroclinic flows through reduced-order stochastic models,” *Commun Math Sci* **16**,
1544 17–51 (2018).
- 1545 ³⁰D. Qi and A. J. Majda, “Predicting fat-tailed intermittent probability distributions in
1546 passive scalar turbulence with imperfect models through empirical information theory,”
1547 *Communications in Mathematical Sciences* **14**, 1687–1722 (2016).
- 1548 ³¹A. J. Majda and D. Qi, “Strategies for reduced-order models for predicting the statistical
1549 responses and uncertainty quantification in complex turbulent dynamical systems,” *SIAM*
1550 *Review* **60**, 491–549 (2018).
- 1551 ³²M. Farazmand and T. Sapsis, “Extreme events: Mechanisms and prediction,” *Applied*
1552 *Mechanics Reviews* (2018).
- 1553 ³³T. P. Sapsis, “New perspectives for the prediction and statistical quantification of extreme
1554 events in high-dimensional dynamical systems,” *Philosophical Transactions of the Royal*
1555 *Society A: Mathematical, Physical and Engineering Sciences* **376**, 20170133 (2018).
- 1556 ³⁴N. Chen and A. Majda, “Conditional Gaussian systems for multiscale nonlinear stochastic
1557 systems: Prediction, state estimation and uncertainty quantification,” *Entropy* **20**, 509
1558 (2018).
- 1559 ³⁵A. J. Majda and J. Harlim, “Physics constrained nonlinear regression models for time
1560 series,” *Nonlinearity* **26**, 201 (2012).
- 1561 ³⁶J. Harlim, A. Mahdi, and A. J. Majda, “An ensemble Kalman filter for statistical es-
1562 timation of physics constrained nonlinear regression models,” *Journal of Computational*
1563 *Physics* **257**, 782–812 (2014).
- 1564 ³⁷R. S. Liptser and A. N. Shiryaev, “Statistics of random processes II: Applications,” *Appl.*

- 1565 Math **6** (2001).
- 1566 ³⁸N. Chen and A. J. Majda, “Filtering nonlinear turbulent dynamical systems through
1567 conditional Gaussian statistics,” *Monthly Weather Review* **144**, 4885–4917 (2016).
- 1568 ³⁹D. Giannakis and A. J. Majda, “Quantifying the predictive skill in long-range forecast-
1569 ing. Part II: Model error in coarse-grained Markov models with application to ocean-
1570 circulation regimes,” *Journal of Climate* **25**, 1814–1826 (2012).
- 1571 ⁴⁰J. S. Armstrong, *Long-range forecasting* (Wiley New York ETC., 1985).
- 1572 ⁴¹H. H. Hendon, B. Liebmann, M. Newman, J. D. Glick, and J. Schemm, “Medium-range
1573 forecast errors associated with active episodes of the Madden–Julian oscillation,” *Monthly
1574 Weather Review* **128**, 69–86 (2000).
- 1575 ⁴²P. R. Vlachas, W. Byeon, Z. Y. Wan, T. P. Sapsis, and P. Koumoutsakos, “Data-driven
1576 forecasting of high-dimensional chaotic systems with long short-term memory networks,”
1577 *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*
1578 **474**, 20170844 (2018).
- 1579 ⁴³Z. Y. Wan and T. P. Sapsis, “Reduced-space Gaussian process regression for data-driven
1580 probabilistic forecast of chaotic dynamical systems,” *Physica D: Nonlinear Phenomena*
1581 **345**, 40–55 (2017).
- 1582 ⁴⁴S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation* **9**,
1583 1735–1780 (1997).
- 1584 ⁴⁵K. P. Murphy, *Machine learning: a probabilistic perspective* (MIT press, 2012).
- 1585 ⁴⁶S. Haykin, *Neural networks: a comprehensive foundation* (Prentice Hall PTR, 1994).
- 1586 ⁴⁷B. Scholkopf and A. J. Smola, *Learning with kernels: support vector machines, regular-
1587 ization, optimization, and beyond* (MIT press, 2001).
- 1588 ⁴⁸B. Gershgorin, J. Harlim, and A. J. Majda, “Improving filtering and prediction of spa-
1589 tially extended turbulent systems with model errors through stochastic parameter esti-
1590 mation,” *Journal of Computational Physics* **229**, 32–57 (2010).
- 1591 ⁴⁹B. Gershgorin, J. Harlim, and A. J. Majda, “Test models for improving filtering with
1592 model errors through stochastic parameter estimation,” *Journal of Computational Physics*
1593 **229**, 1–31 (2010).
- 1594 ⁵⁰M. Branicki and A. J. Majda, “Dynamic stochastic superresolution of sparsely observed
1595 turbulent systems,” *Journal of Computational Physics* **241**, 333–363 (2013).
- 1596 ⁵¹S. R. Keating, A. J. Majda, and K. S. Smith, “New methods for estimating ocean

- 1597 eddy heat transport using satellite altimetry,” *Monthly Weather Review* **140**, 1703–1722
1598 (2012).
- 1599 ⁵²A. J. Majda and I. Grooms, “New perspectives on superparameterization for geophysical
1600 turbulence,” *Journal of Computational Physics* **271**, 60–77 (2014).
- 1601 ⁵³A. J. Majda, I. Timofeyev, and E. V. Eijnden, “Models for stochastic climate prediction,”
1602 *Proceedings of the National Academy of Sciences* **96**, 14687–14691 (1999).
- 1603 ⁵⁴A. J. Majda, I. Timofeyev, and E. Vanden Eijnden, “A mathematical framework for
1604 stochastic climate models,” *Communications on Pure and Applied Mathematics* **54**, 891–
1605 974 (2001).
- 1606 ⁵⁵A. Majda, I. Timofeyev, and E. Vanden-Eijnden, “A priori tests of a stochastic mode
1607 reduction strategy,” *Physica D: Nonlinear Phenomena* **170**, 206–252 (2002).
- 1608 ⁵⁶A. J. Majda, I. Timofeyev, and E. Vanden-Eijnden, “Systematic strategies for stochastic
1609 mode reduction in climate,” *Journal of the Atmospheric Sciences* **60**, 1705–1722 (2003).
- 1610 ⁵⁷A. J. Majda and P. Embid, “Averaging over fast gravity waves for geophysical flows
1611 with unbalanced initial data,” *Theoretical and computational fluid dynamics* **11**, 155–
1612 169 (1998).
- 1613 ⁵⁸N. Chen, A. J. Majda, and D. Giannakis, “Predicting the cloud patterns of the Madden-
1614 Julian oscillation through a low-order nonlinear stochastic model,” *Geophysical Research*
1615 *Letters* **41**, 5612–5619 (2014).
- 1616 ⁵⁹N. Chen and A. J. Majda, “Predicting the real-time multivariate Madden–Julian oscil-
1617 lation index through a low-order nonlinear stochastic model,” *Monthly Weather Review*
1618 **143**, 2148–2169 (2015).
- 1619 ⁶⁰N. Chen and A. J. Majda, “Predicting the cloud patterns for the boreal summer in-
1620 traseasonal oscillation through a low-order stochastic model,” *Mathematics of Climate*
1621 *and Weather Forecasting* **1**, 1–20 (2015).
- 1622 ⁶¹N. Chen and A. J. Majda, “Filtering the stochastic skeleton model for the Madden–Julian
1623 oscillation,” *Monthly Weather Review* **144**, 501–527 (2016).
- 1624 ⁶²N. Chen, A. J. Majda, and X. T. Tong, “Information barriers for noisy Lagrangian tracers
1625 in filtering random incompressible flows,” *Nonlinearity* **27**, 2133 (2014).
- 1626 ⁶³N. Chen, A. J. Majda, and X. T. Tong, “Noisy Lagrangian tracers for filtering random
1627 rotating compressible flows,” *Journal of Nonlinear Science* **25**, 451–488 (2015).
- 1628 ⁶⁴N. Chen and A. J. Majda, “Model error in filtering random compressible flows utilizing

noisy Lagrangian tracers,” *Monthly Weather Review* **144**, 4037–4061 (2016).

⁶⁵A. J. Majda, D. Qi, and T. P. Sapsis, “Blended particle filters for large-dimensional chaotic dynamical systems,” *Proceedings of the National Academy of Sciences*, 201405675 (2014).

⁶⁶N. Chen and A. J. Majda, “Efficient statistically accurate algorithms for the Fokker–Planck equation in large dimensions,” *Journal of Computational Physics* **354**, 242–268 (2018).

⁶⁷N. Chen, A. J. Majda, and X. T. Tong, “Rigorous analysis for efficient statistically accurate algorithms for solving Fokker–Planck equations in large dimensions,” *SIAM/ASA Journal on Uncertainty Quantification* **6**, 1198–1223 (2018).

⁶⁸N. Chen and A. J. Majda, “Beating the curse of dimension with accurate statistics for the Fokker–Planck equation in complex turbulent systems,” *Proceedings of the National Academy of Sciences* **114**, 12864–12869 (2017).

⁶⁹R. E. Kalman, “A new approach to linear filtering and prediction problems,” *Journal of basic Engineering* **82**, 35–45 (1960).

⁷⁰R. E. Kalman and R. S. Bucy, “New results in linear filtering and prediction theory,” *Journal of basic engineering* **83**, 95–108 (1961).

⁷¹K. Brammer and G. Siffing, *Kalman-Bucy Filters* (Artech House on Demand, 1989).

⁷²E. N. Lorenz, “Deterministic nonperiodic flow,” *Journal of the atmospheric sciences* **20**, 130–141 (1963).

⁷³E. N. Lorenz, “Predictability: A problem partly solved,” in *Proc. Seminar on predictability*, Vol. 1 (1996).

⁷⁴A. J. Majda and M. Branicki, “Lessons in uncertainty quantification for turbulent dynamical systems,” *Discrete & Continuous Dynamical Systems-A* **32**, 3133–3221 (2012).

⁷⁵A. J. Majda and B. Gershgorin, “Improving model fidelity and sensitivity for complex systems through empirical information theory,” *Proceedings of the National Academy of Sciences* **108**, 10044–10049 (2011).

⁷⁶J. von Neumann, “Some remarks on the problem of forecasting climatic fluctuations,” in *Dynamics of climate* (Elsevier, 1960) pp. 9–11.

⁷⁷M. W. Moncrieff, M. A. Shapiro, J. M. Slingo, and F. Molteni, “Collaborative research at the intersection of weather and climate,” *Bulletin of the World Meteorological Organization* **56**, 204–211 (2007).

- 1661 ⁷⁸R. A. Madden and P. R. Julian, “Detection of a 40–50 day oscillation in the zonal wind
1662 in the tropical pacific,” *Journal of the atmospheric sciences* **28**, 702–708 (1971).
- 1663 ⁷⁹R. A. Madden and P. R. Julian, “Description of global-scale circulation cells in the tropics
1664 with a 40–50 day period,” *Journal of the atmospheric sciences* **29**, 1109–1123 (1972).
- 1665 ⁸⁰A. J. Majda, S. N. Stechmann, S. Chen, R. H. Ogrosky, and S. Thual, *Tropical Intraseasonal Variability and the Stochastic Skeleton Method* (Springer, 2019).
- 1666
1667 ⁸¹N. Chen and A. J. Majda, “A new efficient parameter estimation algorithm for high-
1668 dimensional complex nonlinear turbulent dynamical systems with partial observations,”
1669 *Journal of Computational Physics* (2019), accepted.
- 1670 ⁸²M. A. Tanner and W. H. Wong, “The calculation of posterior distributions by data aug-
1671 mentation,” *Journal of the American statistical Association* **82**, 528–540 (1987).
- 1672 ⁸³A. Golightly and D. J. Wilkinson, “Bayesian inference for nonlinear multivariate diffusion
1673 models observed with error,” *Computational Statistics & Data Analysis* **52**, 1674–1693
1674 (2008).
- 1675 ⁸⁴G. C. Wei and M. A. Tanner, “A Monte Carlo implementation of the EM algorithm
1676 and the poor man’s data augmentation algorithms,” *Journal of the American statistical
1677 Association* **85**, 699–704 (1990).
- 1678 ⁸⁵O. Papaspiliopoulos, G. O. Roberts, and O. Stramer, “Data augmentation for diffusions,”
1679 *Journal of Computational and Graphical Statistics* **22**, 665–688 (2013).
- 1680 ⁸⁶M. Vihola, “Robust adaptive Metropolis algorithm with coerced acceptance rate,” *Statis-
1681 tics and Computing* **22**, 997–1008 (2012).
- 1682 ⁸⁷K. E. Taylor, “Summarizing multiple aspects of model performance in a single diagram,”
1683 *Journal of Geophysical Research: Atmospheres* **106**, 7183–7192 (2001).
- 1684 ⁸⁸P. L. Houtekamer and H. L. Mitchell, “Data assimilation using an ensemble Kalman filter
1685 technique,” *Monthly Weather Review* **126**, 796–811 (1998).
- 1686 ⁸⁹P. F. Lermusiaux, “Data assimilation via error subspace statistical estimation. Part I-
1687 I: Middle Atlantic Bight shelfbreak front simulations and ESSE validation,” *Monthly
1688 Weather Review* **127**, 1408–1432 (1999).
- 1689 ⁹⁰H. H. Hendon, E. Lim, G. Wang, O. Alves, and D. Hudson, “Prospects for predicting
1690 two flavors of El Niño,” *Geophysical Research Letters* **36** (2009).
- 1691 ⁹¹H.-M. Kim, P. J. Webster, and J. A. Curry, “Seasonal prediction skill of ECMWF System
1692 4 and NCEP CFSv2 retrospective forecast for the northern hemisphere winter,” *Climate*

- 1693 Dynamics **39**, 2957–2973 (2012).
- 1694 ⁹²A. J. Majda and D. Qi, “Linear and nonlinear statistical response theories with prototype
1695 applications to sensitivity analysis and statistical control of complex turbulent dynamical
1696 systems,” *Chaos: An Interdisciplinary Journal of Nonlinear Science* (2019), submitted.
- 1697 ⁹³D. Qi and A. J. Majda, “Low-dimensional reduced-order models for statistical response
1698 and uncertainty quantification: Two-layer baroclinic turbulence,” *Journal of the Atmo-
1699 spheric Sciences* **73**, 4609–4639 (2016).
- 1700 ⁹⁴A. J. Majda and B. Gershgorin, “Quantifying uncertainty in climate change science
1701 through empirical information theory,” *Proceedings of the National Academy of Sciences*
1702 **107**, 14958–14963 (2010).
- 1703 ⁹⁵R. Kleeman, “Measuring dynamical prediction utility using relative entropy,” *Journal of
1704 the atmospheric sciences* **59**, 2057–2072 (2002).
- 1705 ⁹⁶R. Kleeman, “Information theory and dynamical system predictability,” *Entropy* **13**, 612–
1706 649 (2011).
- 1707 ⁹⁷M. Branicki and A. Majda, “Quantifying Bayesian filter performance for turbulent dy-
1708 namical systems through information theory,” *Commun. Math. Sci* **12**, 901–978 (2014).
- 1709 ⁹⁸A. Majda, R. V. Abramov, and M. J. Grote, *Information theory and stochastics for
1710 multiscale nonlinear systems*, Vol. 25 (American Mathematical Soc., 2005).
- 1711 ⁹⁹S. Kullback and R. A. Leibler, “On information and sufficiency,” *The annals of mathe-
1712 matical statistics* **22**, 79–86 (1951).
- 1713 ¹⁰⁰S. Kullback, “Letter to the editor: The Kullback-Leibler distance,” *AMERICAN STATIS-
1714 TICIEN* (1987).
- 1715 ¹⁰¹S. Kullback, “*Statistics and information theory*,” J Wiley Sons, New York (1959).
- 1716 ¹⁰²N. Chen, D. Giannakis, R. Herbei, and A. J. Majda, “An MCMC algorithm for param-
1717 eter estimation in signals with hidden intermittent instability,” *SIAM/ASA Journal on
1718 Uncertainty Quantification* **2**, 647–669 (2014).
- 1719 ¹⁰³D. Giannakis, A. J. Majda, and I. Horenko, “Information theory, model error, and
1720 predictive skill of stochastic models for complex nonlinear systems,” *Physica D: Nonlinear
1721 Phenomena* **241**, 1735–1752 (2012).
- 1722 ¹⁰⁴T. N. Palmer, “A nonlinear dynamical perspective on model error: A proposal for non-
1723 local stochastic-dynamic parametrization in weather and climate prediction models,”
1724 *Quarterly Journal of the Royal Meteorological Society* **127**, 279–304 (2001).

- 1725 ¹⁰⁵J. Berner, U. Achatz, L. Batté, L. Bengtsson, A. d. l. Cámara, H. M. Christensen,
1726 M. Colangeli, D. R. Coleman, D. Crommelin, S. I. Dolaptchiev, *et al.*, “Stochastic param-
1727 eterization: Toward a new view of weather and climate models,” *Bulletin of the American*
1728 *Meteorological Society* **98**, 565–588 (2017).
- 1729 ¹⁰⁶J. W.-B. Lin and J. D. Neelin, “Toward stochastic deep convective parameterization in
1730 general circulation models,” *Geophysical research letters* **30** (2003).
- 1731 ¹⁰⁷M. Branicki and A. J. Majda, “Quantifying uncertainty for predictions with model error
1732 in non-Gaussian systems with intermittency,” *Nonlinearity* **25**, 2543 (2012).
- 1733 ¹⁰⁸A. J. Majda, R. Abramov, and B. Gershgorin, “High skill in low-frequency climate re-
1734 sponse through fluctuation dissipation theorems despite structural instability,” *Proceed-*
1735 *ings of the National Academy of Sciences* **107**, 581–586 (2010).
- 1736 ¹⁰⁹A. J. Majda, C. Franzke, and D. Crommelin, “Normal forms for reduced stochastic
1737 climate models,” *Proceedings of the National Academy of Sciences* **106**, 3649–3653 (2009).
- 1738 ¹¹⁰A. J. Majda and Y. Yuan, “Fundamental limitations of ad hoc linear and quadratic
1739 multi-level regression models for physical systems,” *Discrete and Continuous Dynamical*
1740 *Systems B* **17**, 1333–1363 (2012).
- 1741 ¹¹¹J. L. Anderson, “An ensemble adjustment Kalman filter for data assimilation,” *Monthly*
1742 *weather review* **129**, 2884–2903 (2001).
- 1743 ¹¹²C. Sparrow, *The Lorenz equations: bifurcations, chaos, and strange attractors*, Vol. 41
1744 (Springer Science & Business Media, 2012).
- 1745 ¹¹³H. Haken, “Analogy between higher instabilities in fluids and lasers,” *Physics Letters A*
1746 **53**, 77–78 (1975).
- 1747 ¹¹⁴E. Knobloch, “Chaos in the segmented disc dynamo,” *Physics Letters A* **82**, 439–440
1748 (1981).
- 1749 ¹¹⁵M. Gorman, P. Widmann, and K. Robbins, “Nonlinear dynamics of a convection loop:
1750 a quantitative comparison of experiment with theory,” *Physica D: Nonlinear Phenomena*
1751 **19**, 255–267 (1986).
- 1752 ¹¹⁶N. Hemati, “Strange attractors in brushless DC motors,” *IEEE Transactions on Circuits*
1753 *and Systems I: Fundamental Theory and Applications* **41**, 40–45 (1994).
- 1754 ¹¹⁷K. M. Cuomo and A. V. Oppenheim, “Circuit implementation of synchronized chaos with
1755 applications to communications,” *Physical review letters* **71**, 65 (1993).
- 1756 ¹¹⁸D. Poland, “Cooperative catalysis and chemical chaos: a chemical model for the Lorenz

- 1757 equations,” *Physica D: Nonlinear Phenomena* **65**, 86–99 (1993).
- 1758 ¹¹⁹S. I. Tzenov, “Strange attractors characterizing the osmotic instability,” arXiv preprint
1759 arXiv:1406.0979 (2014).
- 1760 ¹²⁰J. G. Charney and J. G. DeVore, “Multiple flow equilibria in the atmosphere and block-
1761 ing,” *Journal of the atmospheric sciences* **36**, 1205–1216 (1979).
- 1762 ¹²¹A. Wiin-Nielsen, “Steady states and stability properties of a low-order barotropic system
1763 with forcing and dissipation,” *Tellus* **31**, 375–386 (1979).
- 1764 ¹²²C. Franzke, D. Crommelin, A. Fischer, and A. J. Majda, “A hidden Markov model
1765 perspective on regimes and metastability in atmospheric flows,” *Journal of Climate* **21**,
1766 1740–1757 (2008).
- 1767 ¹²³F. Baur, “Extended-range weather forecasting,” in *Compendium of meteorology* (Springer,
1768 1951) pp. 814–833.
- 1769 ¹²⁴D. F. Rex, “Blocking action in the middle troposphere and its effect upon regional climate:
1770 I. An aerological study of blocking action,” *Tellus* **2**, 196–211 (1950).
- 1771 ¹²⁵R. M. Dole and N. D. Gordon, “Persistent anomalies of the extratropical northern hemi-
1772 sphere wintertime circulation: Geographical distribution and regional persistence charac-
1773 teristics,” *Monthly Weather Review* **111**, 1567–1586 (1983).
- 1774 ¹²⁶X. Cheng and J. M. Wallace, “Cluster analysis of the northern hemisphere wintertime 500-
1775 hPa height field: Spatial patterns,” *Journal of the Atmospheric Sciences* **50**, 2674–2696
1776 (1993).
- 1777 ¹²⁷M. Kimoto and M. Ghil, “Multiple flow regimes in the northern hemisphere winter. part
1778 II: Sectorial regimes and preferred transitions,” *Journal of the atmospheric sciences* **50**,
1779 2645–2673 (1993).
- 1780 ¹²⁸C. Franzke and A. J. Majda, “Low-order stochastic mode reduction for a prototype at-
1781 mospheric GCM,” *Journal of the atmospheric sciences* **63**, 457–479 (2006).
- 1782 ¹²⁹J. Berner and G. Branstator, “Linear and nonlinear signatures in the planetary wave dy-
1783 namics of an AGCM: Probability density functions,” *Journal of the atmospheric sciences*
1784 **64**, 117–136 (2007).
- 1785 ¹³⁰A. J. Majda, C. Franzke, and B. Khouider, “An applied mathematics perspective on
1786 stochastic modelling for climate,” *Philosophical Transactions of the Royal Society A:*
1787 *Mathematical, Physical and Engineering Sciences* **366**, 2427–2453 (2008).
- 1788 ¹³¹D. Qi and A. J. Majda, “Low-dimensional reduced-order models for statistical response

- 1789 and uncertainty quantification: Barotropic turbulence with topography,” *Physica D: Non-*
1790 *linear Phenomena* **343**, 7–27 (2017).
- 1791 ¹³²D. Qi and A. J. Majda, “Rigorous statistical bounds in uncertainty quantification for one-
1792 layer turbulent geophysical flows,” *Journal of Nonlinear Science* **28**, 1709–1761 (2018).
- 1793 ¹³³D. T. Crommelin, J. Opsteegh, and F. Verhulst, “A mechanism for atmospheric regime
1794 behavior,” *Journal of the atmospheric sciences* **61**, 1406–1419 (2004).
- 1795 ¹³⁴D. Crommelin and A. Majda, “Strategies for model reduction: Comparing different opti-
1796 mal bases,” *Journal of the Atmospheric Sciences* **61**, 2206–2217 (2004).