

## Research Articles

1

Nan Chen\* and Andrew J Majda

2

# Predicting the Cloud Patterns for the Boreal Summer Intraseasonal Oscillation Through a Low-Order Stochastic Model

3

4

5

**Abstract:** We assess the predictability limits of the large-scale cloud patterns in the boreal summer intraseasonal variability (BSISO), which are measured by the infrared brightness temperature, a proxy for convective activity. A recent developed nonlinear data analysis technique, nonlinear Laplacian spectrum analysis (NLSA), is applied to the brightness temperature data, defining two spatial modes with high intermittency associated with the BSISO time series. Then a recent developed data-driven physics-constrained low-order modeling strategy is applied to these time series. The result is a four dimensional system with two observed BSISO variables and two hidden variables involving correlated multiplicative noise through the nonlinear energy-conserving interaction. With the optimal parameters calibrated by information theory, the non-Gaussian fat tailed probability distribution functions (PDFs), the autocorrelations and the power spectrum of the model signals almost perfectly match those of the observed data. An ensemble prediction scheme incorporating an effective on-line data assimilation algorithm for determining the initial ensemble of the hidden variables shows the useful prediction skill in the non-El Niño years is at least 30 days and even reaches 55 days in those years with regular oscillations and the skillful prediction lasts for 18 days in the strong El Niño year (year 1998). Furthermore, the ensemble spread succeeds in indicating the forecast uncertainty. Although the reduced linear model with time-periodic stable-unstable damping is able to capture the non-Gaussian fat tailed PDFs, it is less skillful in forecasting the BSISO in the years with irregular oscillations. The failure of the ensemble spread to include the truth also indicates failure in quantification of the uncertainty. In addition, without the energy-conserving nonlinear interactions, the linear model is sensitive with parameter variations.

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

---

**\*Corresponding author: Nan Chen:** Department of Mathematics and Center for Atmosphere Ocean Science, Courant Institute of Mathematical Sciences, New York University, New York, USA. [chennan@cims.nyu.edu](mailto:chennan@cims.nyu.edu)

**Andrew J Majda:** Department of Mathematics and Center for Atmosphere Ocean Science, Courant Institute of Mathematical Sciences, New York University, New York, USA and Center for Prototype Climate Modeling, NYU Abu Dhabi, Saadiyat Island, Abu Dhabi, UAE. [jonjon@cims.nyu.edu](mailto:jonjon@cims.nyu.edu)

1 Finally, the twin experiment with nonlinear stochastic model has comparable  
2 skill as the observed data, suggesting the nonlinear stochastic model has  
3 significant skill for determining the predictability limits of the large-scale cloud  
4 patterns of the BSISO.

5 **Keywords:** Nonlinear Laplacian spectrum analysis (NLSA), data-driven physics-  
6 constrained low-order modeling strategy, information theory, data assimilation,  
7 predictability limits

8 **PACS:** 60G25, 60H10, 76E20, 94A15

9 DOI: ..., Received ...; accepted ...

10 —————

11 **Communicated by:** ...

12 **Für** ...

## 13 1 Introduction

14 The boreal summer intraseasonal oscillation (BSISO) is one of the prominent  
15 modes of tropical intraseasonal variability. As a slow moving planetary scale  
16 envelope of convection propagating northward [1, 2], the BSISO distinguishes  
17 itself from the Madden-Julian oscillation (MJO), which prevails during boreal  
18 winter and propagates eastward. The BSISO is known to affect summer  
19 monsoon onset and active/break phases [3, 4, 5], and the seasonal means  
20 of summer monsoons [6, 7]. It also has fundamental impacts on the tropical  
21 precipitation, the frequency of tropical cyclones and extra-tropical climate  
22 variations [8, 9]. The studies of the prediction of the BSISO are mostly through  
23 operational dynamical models [1, 10, 11, 12, 13, 14, 15, 16] with only a few  
24 works focusing on low-order statistical models [17, 18, 19, 20]. One reason for  
25 the lack of utilizing low-order models for prediction is that different from the  
26 MJO (e.g., real-time multivariate MJO (RMM) index [21]) low-dimensional  
27 real-time monitoring and forecast verification metrics for the BSISO were not  
28 available until the recent time (since 2013) [22, 23]. The BSISO index in [23] is  
29 based on an extended empirical orthogonal function (EEOF) analysis on daily  
30 unfiltered rainfall anomalies. The derivation of the BSISO index in [22] mimics  
31 that for the RMM index and is based on the multivariate EOF analysis of  
32 daily anomalies of the zonal wind at 850 hPa and outgoing long-wave radiation  
33 (OLR), a proxy for convective activity. While the use these indices improves  
34 the quantification of skill of extended range forecasts of the BSISO, these linear

data analysis methods have limitations in capturing the highly nonlinear and 1  
intermittent characters of the BSISO [24]. 2

Here we assess the predictability limits of the BSISO of the large-scale 3  
cloud patterns, which are measured by the infrared brightness temperature, 4  
a proxy for convective activity, alone. This is achieved in two steps. In the first 5  
step, a recent developed advanced nonlinear data analysis technique, nonlinear 6  
Laplacian spectrum analysis (NLSA), is applied to the brightness temperature 7  
data to define two spatial patterns associated with the BSISO. A key advantage 8  
of NLSA is that it requires no preprocessing such as bandpass filtering or seasonal 9  
partitioning of the input data, enabling simultaneous recovery of the dominant 10  
BSISO modes [25, 26, 27, 28, 29]. NLSA provides two time series representing 11  
the BSISO. These two BSISO time series are highly intermittent with non- 12  
Gaussian fat tailed probability distribution function (PDF), which differs from 13  
those derived by straightforward linear methods [22, 23]. The second step is to 14  
apply a recent systematic strategy for data-driven physics-constrained low-order 15  
stochastic modeling of time series to the two BSISO time series [30, 31]. The 16  
result is a four-dimensional nonlinear stochastic model with two observed state 17  
variables representing the two BSISO indices and two hidden variables. This 18  
low-order model involves correlated multiplicative noise through the energy- 19  
conserving nonlinear interaction between the observed and hidden variables as 20  
well as the additive noise. Note that this nonlinear low-order stochastic model 21  
has been shown to have significant skill for determining the predictability limits 22  
of the large-scale cloud patterns of the boreal winter MJO [32]. In addition, 23  
incorporating a new information-theoretic strategy in the training phase, a 24  
slight simplified version of (1) has been adopted to significantly improve the 25  
predictability of the real-time multivariate MJO indices [33]. 26

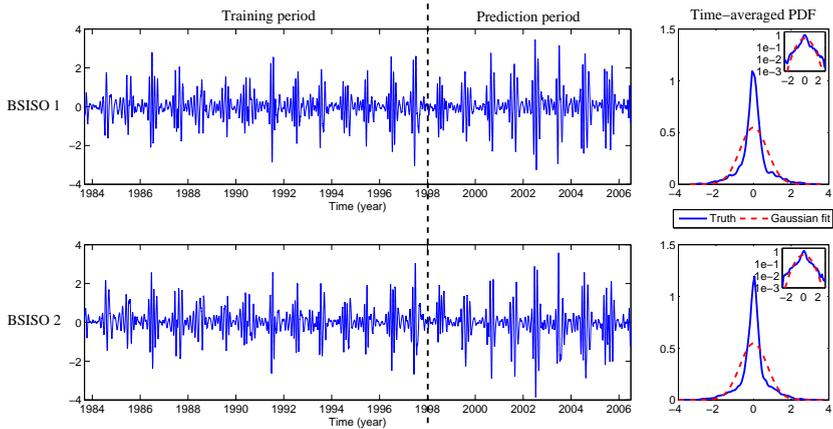
The remainder of the paper is organized as follows. Section 2 describes the 27  
source of the brightness temperature dataset and the application of NLSA to 28  
the brightness temperature data to form the two BSISO time series. Section 29  
3 involves the nonlinear low-order stochastic model as well as the information- 30  
theoretic calibration strategy and data assimilation algorithm for determining 31  
the initial ensemble of the hidden variables in the ensemble prediction scheme. 32  
This is followed by the prediction results, shown in Section 4. To understand 33  
the role of the nonlinearity in the low-order model, the linear model as well as 34  
its forecasting skill is studied in Section 5. Section 6 includes the perfect model 35  
twin experiment in checking the model error and the predictability limits. The 36  
paper is concluded in Section 7. 37

## 1 2 The Boreal Summer Intraseasonal Oscillation 2 Through NLSA

3 Here tropical convection is analyzed in satellite observations of infrared brightness  
4 temperature ( $T_b$ ) from the Cloud Archive User Service (CLAUS) Version 4.7  
5 [34], which is sampled every 3 hours from 1 July 1983 to 30 June 2006 with a  
6 spatial resolution of  $0.5^\circ$ . Brightness temperature is a measure of the earth's  
7 infrared emission in terms of the temperature of a hypothesized blackbody  
8 emitting the same amount of radiation at the same wavelength ( $\sim 10\text{--}11 \mu\text{m}$  in  
9 CLAUS). It is a highly correlated variable with the total terrestrial longwave  
10 emission. In the tropics, positive (negative)  $T_b$  anomalies are associated with  
11 reduced (increased) cloudiness, hence suppressed (enhanced) deep convection.

12 The NLSA algorithm is applied to the full CLAUS dataset restricted to the  
13 tropical belt  $15^\circ\text{N}\text{--}15^\circ\text{S}$ , with a lagged embedding window of 60 days. A variety  
14 of extended spatial cloud patterns emerge from the analysis, including annual,  
15 interannual, intraseasonal and diurnal modes, but the focus here is on the two  
16 spatial cloud patterns with time series depicted in Figure 1. It is clear that  
17 these time series are active from May to September of each year corresponding  
18 to boreal summer. Hereafter, we utilize the terminology, BSISO indices, for the  
19 two time series in Figure 1. The details of the NLSA method is available in  
20 [25, 26, 27, 28, 29]. Applying NLSA to the CLAUS dataset to obtain the BSISO  
21 indices are well documented in [24], where the spatial patterns associated with  
22 the BSISO indices in different phases are shown in Figure 9b in [24] and the  
23 spatiotemporal evolution of BSISO is shown in Movie 1h there. The figure and  
24 movie indicate the character of BSISO that a cluster of positive  $T_b$  develops in  
25 the central Indian Ocean and then moves northward towards the Bay of Bengal  
26 and India and branches off towards the western Pacific and the Monsoon Trough,  
27 bypassing the Maritime Continent from the north. Then following the dry phase  
28 of BSISO, a cluster of anomalously high convection develops in the central Indian  
29 Ocean, and propagates towards India and the western Pacific, completing the  
30 BSISO cycle. In addition, Figure 11 and Movie 2 in [24] compares the BSISO  
31 indices derived from NLSA with those from the EOF-singular spectrum analysis  
32 which is a linear data analysis method. The reconstructed BSISO and MJO  
33 patterns from the EOF method have coarser structures than their NLSA-based  
34 counterparts and appear to mix BSISO and MJO propagation. The highly  
35 non-Gaussian fat-tailed distributions in the BSISO indices from NLSA, which  
36 corresponding to intermittency and strong seasonality, also contrast with the  
37 nearly Gaussian distributions in the EOF related indices. We point out that  
38 the BSISO indices studied in this work is slightly cleaner than those in [24] by

reconstructing the eigenfunctions associated with NLSA through the shift map 1  
 [35] while the corresponding spatiotemporal patterns have almost no changes. 2  
 See Appendix A for details. 3



**Fig. 1.** (Left) BSISO indices from NLSA ranging from 3 September 1983 to 30 June 2006. The time series from September 1983 to December 1997 is utilized as training period to get the statistics and that from January 1998 to December 2005 represents the predicting period which will be predicted by the nonlinear low-order stochastic model (1). (Right) The associated time-averaged probability distribution function (PDF) of each index and its Gaussian fit. Here the time-averaged PDF means we take all the points in each one-dimensional time series and compute the PDF based on these points. The small panel inside each subplot shows the PDF in the logarithm scale.

Note that the derivation of the BSISO indices from NLSA shown in Figure 4  
 1 is based on the dataset within the band  $15^{\circ}\text{N}$ – $15^{\circ}\text{S}$ . Yet, the BSISO activities 5  
 typically propagate northward beyond  $15^{\circ}\text{N}$  [1]. Actually, Movie (1h) in [24] 6  
 also shows that the cluster of convection moves outside the north boundary. To 7  
 check the cluster propagation behavior in a farther northward region, we have 8  
 looked at the projected spatiotemporal patterns in the extended band  $30^{\circ}\text{N}$ – 9  
 $15^{\circ}\text{S}$  with the given BSISO indices in Figure 1. The northward propagation of 10  
 the cluster extending up to  $20^{\circ}\text{N}$ – $30^{\circ}\text{N}$  in the spatiotemporal patterns provides 11  
 the evidence that the time series in Figure 1 are suitable BSISO indices. The 12  
 extension of the NLSA and of this work to a wider band, e.g.,  $30^{\circ}\text{N}$ – $15^{\circ}\text{S}$ , for a 13  
 better coverage of the Northward propagating signals is more involved and it is 14  
 thus left for future investigations. 15

### 1 3 The Nonlinear Physics-Constrained 2 Low-Order Stochastic Model

Denote by  $u_1$  and  $u_2$  the two components, BSISO 1 and BSISO 2, described in Figure 1. The PDFs for  $u_1$  and  $u_2$  are highly non-Gaussian with fat tails indicative of the temporal intermittency in the large-scale cloud patterns associated with the BSISO. To describe the variability of the time series  $u_1$  and  $u_2$ , we propose the following family of low-order stochastic models:

$$\frac{du_1}{dt} = (-d_u u_1 + \gamma(v + v_f(t))u_1 - (a + \omega_u)u_2) + \sigma_u \dot{W}_{u_1}, \quad (1a)$$

$$\frac{du_2}{dt} = (-d_u u_2 + \gamma(v + v_f(t))u_2 + (a + \omega_u)u_1) + \sigma_u \dot{W}_{u_2}, \quad (1b)$$

$$\frac{dv}{dt} = (-d_v v - \gamma(u_1^2 + u_2^2)) + \sigma_v \dot{W}_v, \quad (1c)$$

$$\frac{d\omega_u}{dt} = (-d_\omega \omega_u) + \sigma_\omega \dot{W}_\omega, \quad (1d)$$

3 where

$$v_f(t) = f_0 + f_t \sin(\omega_f t + \phi). \quad (2)$$

4 Besides the two observed BSISO variables  $u_1$  and  $u_2$ , the other two variables  
5  $v$  and  $\omega_u$  are hidden and unobserved, representing the stochastic damping and  
6 stochastic phase, respectively. In (1),  $\dot{W}_{u_1}$ ,  $\dot{W}_{u_2}$ ,  $\dot{W}_v$  and  $\dot{W}_\omega$  are independent  
7 white noise. The constant coefficients  $d_u$ ,  $d_v$ , and  $d_\omega$  represent damping for  
8 each stochastic process and have physical dimension  $t^{-1}$ ;  $a$  (also of dimension  
9  $t^{-1}$ ) is the background state of the phases of  $u_1$  and  $u_2$ ;  $\sigma_u$ ,  $\sigma_v$ , and  $\sigma_\omega$  are  
10 noise amplitudes with dimension  $t^{-1/2}$ ; the non-dimensional constant  $\gamma$  is the  
11 coefficient of the nonlinear interaction. The time periodic damping  $v_f(t)$  in  
12 the equations in (1a) and (1b) is utilized to crudely model the active phase  
13 of the BSISO and the quiescent winter season in the seasonal cycle. The  
14 constant coefficients  $\omega_f$  and  $\phi$  in (2) are the frequency and phase of the damping,  
15 respectively. All of the model variables are real.

The hidden variables  $v, \omega_u$  interact with the observed BSISO variables  $u_1, u_2$  through energy-conserving nonlinear interactions following the systematic physics-constrained nonlinear regression strategies for time series developed recently [30, 31]. The energy conserving nonlinear interactions between  $u_1, u_2$  and  $v, \omega_u$  are seen in the following way. First, by dropping the linear and external forcing terms in (1), the remaining equations involving only the nonlinear parts

of (1) read,

$$\frac{du_1}{dt} = \gamma v u_1 - \omega_u u_2, \quad (3a)$$

$$\frac{du_2}{dt} = \gamma v u_2 + \omega_u u_1, \quad (3b)$$

$$\frac{dv}{dt} = -\gamma(u_1^2 + u_2^2), \quad (3c)$$

$$\frac{d\omega_u}{dt} = 0. \quad (3d)$$

To form the evolution equation of the energy from nonlinear interactions  $\tilde{E} = 1$   $(u_1^2 + u_2^2 + v^2 + \omega_u^2)/2$ , we multiply the four equations in (3) by  $u_1, u_2, v$  and  $\omega_u$  2 respectively and then sum them up. The resulting equation yields 3

$$\frac{d\tilde{E}}{dt} = 0. \quad (4)$$

The vanishing of the right hand side in (4) is due to the opposite signs of the 4 nonlinear terms involving  $v$  multiplying  $u_1$  and  $u_2$  in (3a) and (3b) and those in 5 (3c) multiplying by  $v$  as well as the trivial cancelation of skew-symmetric terms 6 involving  $\omega_u$  in (3a) and (3b). 7

The nonlinear low-order stochastic models in (1) are fundamentally different 8 from those utilized earlier [36, 37] in predicting various kinds of climate 9 phenomena which allow for nonlinear interactions only between the observed 10 variables  $u_1, u_2$  and only special linear interactions with layers of hidden 11 variables. The stochastic damping  $v$  and stochastic phase  $\omega_u$  as well as their 12 energy conserving nonlinear interaction with  $u_1$  and  $u_2$  also distinguish the 13 models in (1) from the classic damped harmonic oscillator with only constant 14 damping  $d_u$  and phase  $a$ . It is evident that a negative value of  $\gamma(v + v_f)$  serves 15 to strengthen the total damping of the oscillator. On the other hand, when 16  $\gamma(v + v_f)$  becomes positive and overwhelms  $d_u$ , an exponential growth of  $u_1$  17 and  $u_2$  will occur, which corresponds to the intermittent instability. 18

The nonlinear low-order stochastic model (1) has been shown to have 19 significant skill for determining the predictability limits of the large-scale cloud 20 patterns of the boreal winter MJO [32]. In addition, incorporating a new 21 information-theoretic strategy in the training phase, a simplified version of 22 (1) without the time-period damping  $v_f(t)$  has been adopted to improve the 23 predictability of the real-time multivariate MJO indices [33]. Note that these 24 models are a special case of the models described in [30, 31]. 25

### 1 3.1 Calibration of the nonlinear low-order stochastic models

2 As shown in Figure 1, the full time series are divided into the training (Year  
3 1983-1997) and prediction (Year 1998-2005) periods.

4 The parameters of the stochastic model in (1)–(2) are calibrated by  
5 systematically minimizing the information distance of the highly non-Gaussian  
6 PDFs of the stochastic model compared with that of the actual data [38, 39] and  
7 taking into consideration of the autocorrelations of the two BSISO variables  
8  $u_1, u_2$ . Details are presented in Appendix B which also demonstrates the  
9 robustness of these optimal parameters to their variation. Table 1 records the  
10 optimal parameter values while Figure 2 displays the skill of the stochastic  
11 model with these parameters in recovering the statistics of the two BSISO  
12 indices. Panels (a) and (b) show that the stochastic model (1) succeeds in  
13 capturing the autocorrelations almost perfectly for a three-month duration and  
14 even the wiggles that appears with lags around one year. Panel (c) shows  
15 that the stochastic model captures the fat tailed highly non-Gaussian PDFs of  
16 the two BSISO indices due to intermittency. Panel (d) shows that the power  
17 spectrums of the two BSISO indices from the data and those from the stochastic  
model match very well.

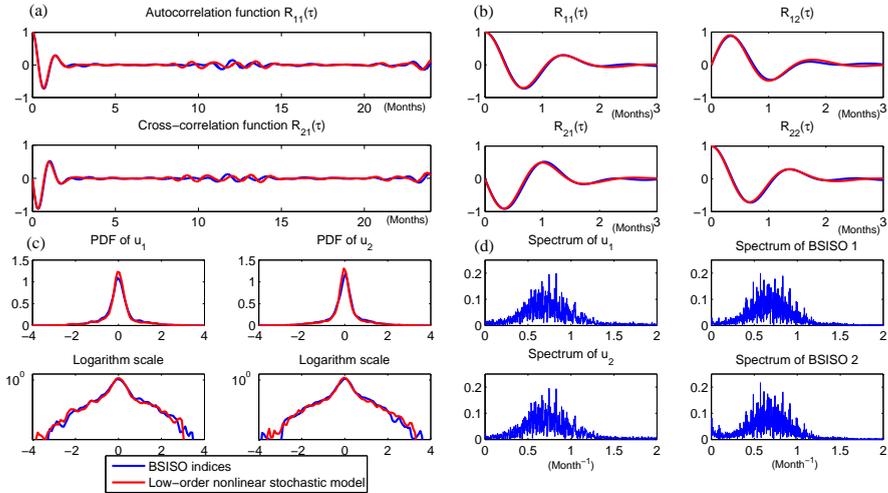
	$d_u$	$d_v$	$d_\omega$	$\sigma_u$	$\sigma_v$	$\sigma_\omega$	$\gamma$	$a$	$f_0$	$f_t$	$\omega_f$	$\phi$
Nonlinear model (1)	0.9	0.9	0.5	0.3	0.8	1.0	0.3	-4.25	1.0	4.0	$2\pi/12$	-3.4
Linear model (5)	0.9			0.35			0.3	-4.25	0.0	4.5	$2\pi/12$	-3.4

**Table 1.** Parameters for the nonlinear low-order stochastic model (1) and linear stochastic model (5). The parameters  $d_u, a, \gamma, f_0, f_t, \omega_f, d_v$  and  $d_\omega$  have units  $m^{-1}$ ;  $\sigma_u, \sigma_v$  and  $\sigma_\omega$  have units  $m^{-1/2}$ ;  $\phi$  is dimensionless. Here  $m$  stands for month.

18

### 19 3.2 Prediction algorithm and data assimilation for the 20 hidden variables

21 The ensemble prediction algorithm is adopted to study the predictability of the  
22 nonlinear low-order stochastic model (1), which involves running the forecast  
23 model (1) forward in time given the initial values. The initial data of the two  
24 state variables  $\mathbf{U} = (u_1, u_2)$  are obtained directly from the observations, i.e.,  
25 BSISO 1 and BSISO 2 indices, and all the ensembles have the same initial values  
26 of  $\mathbf{U} = (u_1, u_2)$ . The more important and challenging issue is to determine



**Fig. 2.** Statistics of the nonlinear low-order stochastic model (1) with the optimal parameters in Table 1. (a) Long-term autocorrelation function  $R_{11}(\tau)$  and cross-correlation function  $R_{21}(\tau)$  from 0 to 24 months. (b) Short-term autocorrelation functions  $R_{11}(\tau)$ ,  $R_{22}(\tau)$  and cross-correlation functions  $R_{21}(\tau)$ ,  $R_{12}(\tau)$  from 0 to 3 months. (c) Equilibrium PDFs of the signal  $u_1, u_2$  from stochastic model compared with that of the BSISO indices. (d) Spectrum of  $u_1, u_2$  compared with that of the BSISO indices.

the initial ensemble of the two hidden variables  $\Gamma = (v, \omega_u)$ . To this end, an 1  
 active data assimilation algorithm is incorporated into the ensemble forecasting 2  
 scheme. 3

The estimates of the hidden parameters  $\Gamma = (v, \omega_u)$  during the training 4  
 period and initialization of these parameters during the prediction phase 5  
 exploit the special structure of the nonlinear low-order stochastic model 6  
 (1). The equations in (1) are a conditional Gaussian system with respect 7  
 to the observations  $\mathbf{U} = (u_1, u_2)$ , meaning that once  $u_1$  and  $u_2$  are given 8  
 the time evolution of the distributions of  $\Gamma = (v, \omega_u)$  is Gaussian. Such 9  
 special feature of (1) allows the closed analytic equations for the conditional 10  
 Gaussian distributions of the hidden parameters  $\Gamma = (v, \omega_u)$  obtained from the 11  
 posterior estimations in the Bayesian framework [40]. Appendix C contains 12  
 the details and explicit equations. We utilize this fact to construct an initial 13  
 ensemble for forecasting at each time in the training and prediction phases for 14  
 $t \in [t_0, t_1, \dots, t_s]$  in the following way. 15

1. Starting from a “burn in” time  $t_-$  earlier than  $t_0$  with arbitrary initial 16  
 conditions for  $\Gamma$ , solve the associated analytic formula (12) until time  $t_0$  17

- 1 to obtain the conditional Gaussian distribution  $p_0(\Gamma|\mathbf{u}(t_0))$ . The initial  
 2 ensemble of the hidden variables  $\Gamma = (v, \omega_u)$  for prediction starting from  $t_0$   
 3 is drawn from this distribution.
- 4 2. The initial ensemble for prediction starting from the next time  $t_1$  is drawn  
 5 from  $p_1(\Gamma|\mathbf{u}(t_1))$ , where  $p_1(\Gamma|\mathbf{u}(t_1))$  is solved by running the analytic  
 6 formula (12) forward from time  $t_0$  to  $t_1$  with the initial value  $p_0(\Gamma|\mathbf{u}(t_0))$ .
- 7 3. Following the same procedure, the initial distributions of the hidden  
 8 variables  $\Gamma = (v, \omega_u)$  for prediction starting from each time  $t_i$  are obtained  
 9 “on the fly” by running the analytic formula (12) forward from time  $t_{i-1}$  to  
 10  $t_i$  with the initial value  $p_{i-1}(\Gamma|\mathbf{u}(t_{i-1}))$  when the new observations up to  
 11  $\mathbf{u}(t_i)$  are available.

12 This is an effective and practical on-line data assimilation algorithm for the  
 13 stochastic models in (1). Note that this algorithm is an improved version of the  
 14 one in [32] for predicting the cloud patterns of the boreal winter MJO, in which  
 15 the initial ensemble at the current time is the combination of all the ensembles  
 16 corresponding to the analogous observations from the historic data and thus the  
 17 algorithm there is more expensive. In the prediction below with (1), we use  $N$   
 18 ensemble members with  $N = 50$ .

## 19 4 Prediction Results

20 With the optimal parameters shown in Table 1 and the effective data assimilation  
 21 algorithm for the ensemble initialization described in Section 3.2, we now study  
 22 the prediction skill of the nonlinear low-order stochastic model (1) in the  
 23 prediction period from year 1998 to year 2005. The BSISO 1 and BSISO 2  
 24 indices from year 1998 to 2005 are shown in panel (a) and (b) of Figure 3. The  
 25 BSISO activities in year 1998 and 1999 are weaker compared with those from  
 26 year 2000 to 2005. Note that year 1998 accompanies strong El Niño events.

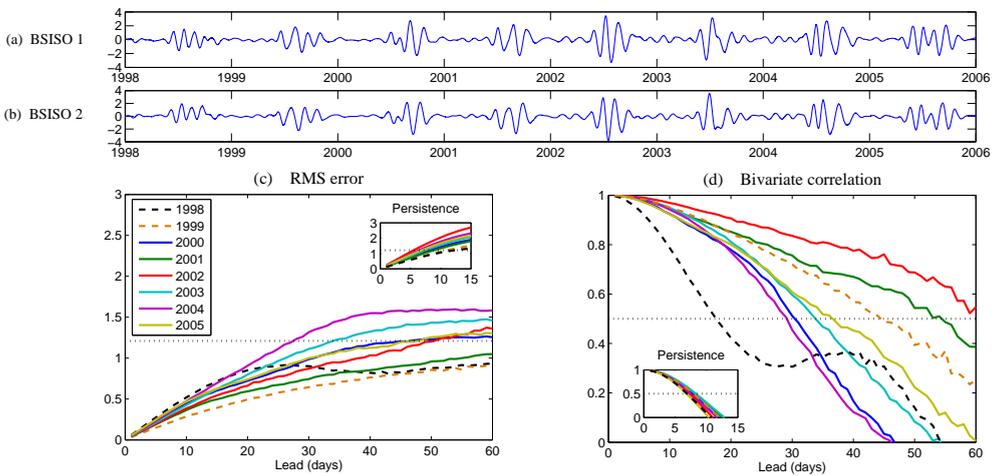
27 Panel (c) and (d) in Figure 3 illustrate the root-mean-squared (RMS) error  
 28 and bivariate correlation in prediction as a function of lead time in different  
 29 years. The threshold level for skillful prediction is indicated by the horizontal  
 30 dotted line in both panels. It is clear that the skillful prediction lasts for more  
 31 than 30 days from year 1999 to 2005. Particularly, the prediction skill in year  
 32 2001 and 2002 with relatively regular intraseasonal oscillations reaches about 55  
 33 days. However, the useful prediction skill in year 1998, 18 days, is significantly  
 34 lower than those in other years. Yet, it is still much improved compared with

the result from persistence prediction, which is at most 8 days for all the years 1  
as shown in the sub-panels in panel (c) and (d). 2

To understand the disparity in the prediction skill in different years, the 3  
predictions of BSISO 2 index regarding the ensemble mean at lead times of 15 4  
and 25 days and 35 and 45 days are shown in Figure 4 and 5, respectively. The 5  
15-day predictions in all years and 25-day predictions in year 1999 to 2005 are 6  
all very skillful ( $\text{Corr} > 0.6$ ), capturing not only the patterns and amplitudes of 7  
the intraseasonal oscillations but some small-scale fluctuations as well. The 8  
prediction skill decreases as the lead time increases as expected. Yet, the 9  
predictions in year 2001 and 2002 at a lead time of 35 days are still quite accurate 10  
( $\text{Corr} > 0.7$ ) and even the 45-day predictions illustrate the similar patterns 11  
( $\text{Corr} > 0.6$ ) as the truth. The unskillful 25-day prediction ( $\text{Corr} = 0.32803$ ) 12  
in year 1998 ascribes to the anomalous feature in BSISO 2 index, that is, the 13  
mean state of the index is significantly above zero (See also Figure 1 for a clearer 14  
visual). Yet, the mean state of the nonlinear oscillator (1), in consistency with 15  
the indices in other years, is around zero. This indicates an intrinsic model 16  
error in describing the BSISO events in year 1998 and in turn results in a low 17  
prediction skill. Actually, strong El Niño events occur at the same the period 18  
as the anomaly in BSISO 2 index. As shown in [41], many examples reveal that 19  
strong El Niño events are coincident with severely weakened BSISO events. The 20  
decrease in the amplitude of BSISO signal implies the increase of the noise to 21  
signal ratio, which prevents the skillful prediction and is consistent with the 22  
findings in [16]. The prediction skill of BSISO 1 is almost the same as that of 23  
BSISO 2, except that BSISO 2 involves a clear demonstration for the intrinsic 24  
model error for predicting the strong El Niño year 1998. In addition, as indicated 25  
in [33], at the strong El Niño phases another intraseasonal oscillation – MJO – 26  
also shows the irregular behaviors and are hard to predict. 27

Figure 6 shows the predictions including the ensemble members starting 28  
from three different dates. The date, April 1, is a time at the transition between 29  
the quiescent phase and the active phase of the BSISO; June 1 is a starting date 30  
in the active mature phase while September 1 is a starting date in the decaying 31  
phase of BSISO activity. Although the ensemble mean forecasts starting from 32  
the April 1 have no long-term skill, which is consistent with the results in Figure 33  
5, the ensemble spread automatically indicates such lack of skill. The envelope 34  
of the ensemble predictions includes the truth for all the years and is a good 35  
indicator of the active and quiescent phases of the BSISO. The ensemble mean 36  
predictions for the June 1 starting date are skillful in short and medium ranges 37  
and the ensemble spreads are the successful uncertainty indicator at long lead 38  
times for all the years. The forecasts starting from September 1 have both an 39  
accurate mean and a small ensemble spread for very long lead times. Note that 40

1 April 1, June 1 and September 1 are the three typical starting days that reflect  
 2 different prediction behaviors. In general, starting from the quiescent phases,  
 3 the ensemble spread accurately tells the active and quiescent phases of BSISO  
 4 although the ensemble mean has no long-term prediction skill for the active  
 5 phase. On the other hand, starting from the active phases, the short and medium  
 6 range forecasting skill is obtained by the ensemble mean and the uncertainty at  
 7 the long lead time is well indicated by the ensemble spread. Therefore, these  
 8 results implies that the important target of predicting the onset and demise time  
 9 of BSISO is well predicted by the ensemble spread.



**Fig. 3.** (Panel (a) and (b)) BSISO indices in the prediction period. (Panel (c) and (d)) Skill scores with RMS error and bivariate correlation for ensemble mean prediction as a function of lead days in different years. The sub-panels inside panel (c) and (d) show the persistence prediction skill up to 15 lead days. The black dotted horizontal line in panel (c) shows the standard deviation of the BSISO indices, i.e., the climatological forecast skill, and that in panel (d) indicates  $\text{Corr} = 0.5$ , the typical threshold for skillful prediction.

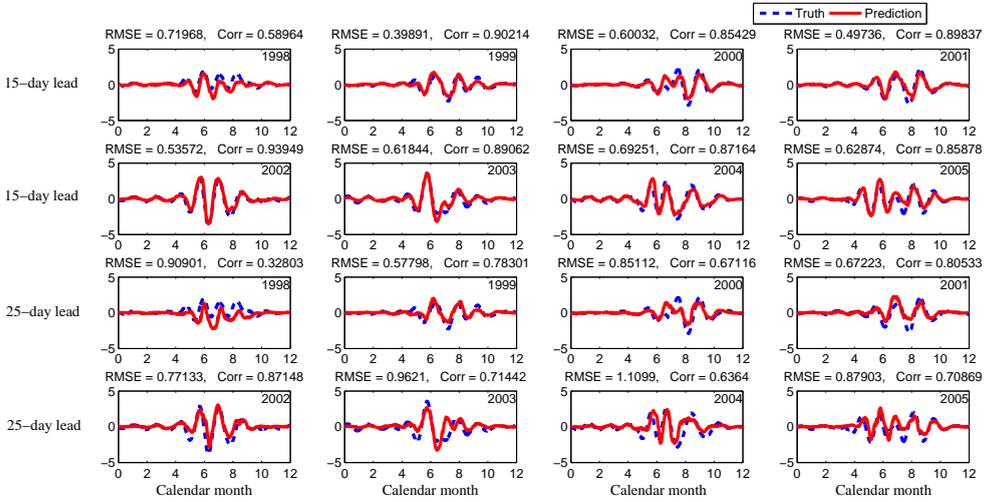


Fig. 4. Prediction of BSISO 2 index regarding the ensemble mean at lead times of 15 and 25 days in different years. The skill scores are computed from both BSISO 1 and 2 indices. The blue curve shows the truth and the red curve shows the ensemble average of the predicted signal. The number of ensemble utilized for prediction is 50.

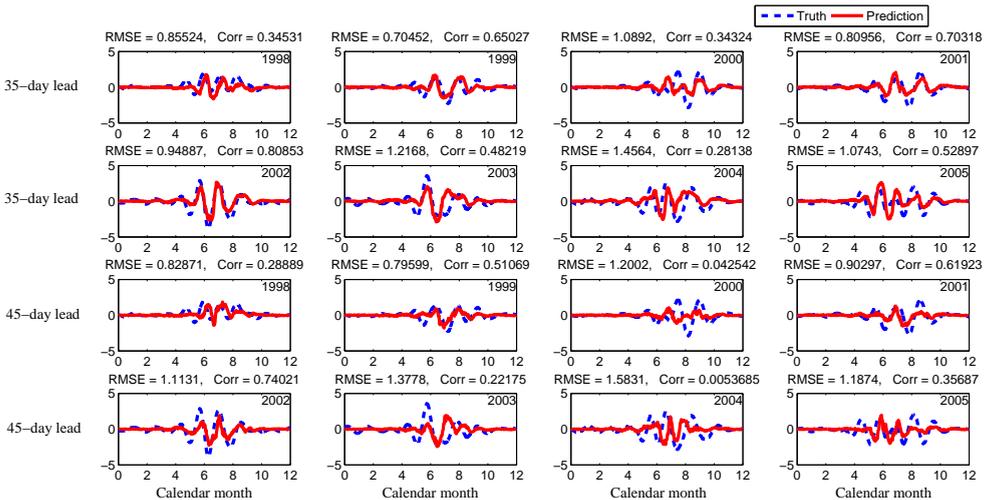
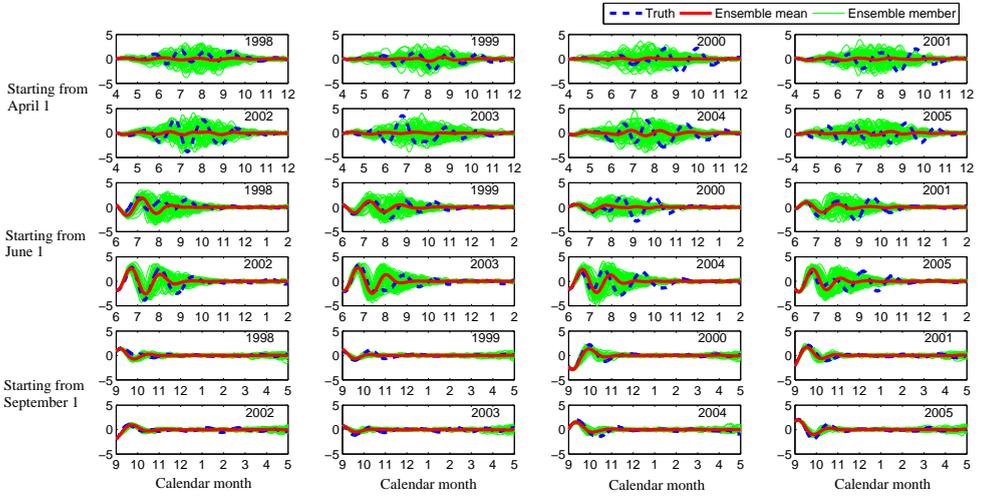


Fig. 5. Same as 4 but at lead times of 35 and 45 days.



**Fig. 6.** (First and second rows) Prediction of BSISO 2 starting from April 1 for different years. Each panel shows the prediction skill of 8 months with the label in x axis indicating the calendar month. (Third and fourth rows) Same as first and second rows but starting from June 1. (Fifth and sixth rows) Same as first and second rows but starting from September 1. The thick blue dashed curve is the BSISO 2 index. The thick red solid curve is the ensemble mean with 50 members, which are shown by the thin solid green curves.

## 5 The Role of the Nonlinearity in the Low-Order Models

1  
2

To understand the role of the nonlinearity in the nonlinear low-order stochastic model (1), we study the prediction skill of the reduced linear version of (1) by removing the stochastic damping  $v$  and stochastic phase  $\omega_u$ ,

$$\frac{du_1}{dt} = (-d_u u_1 + \gamma v_f(t) u_1 - a u_2) + \sigma_u \dot{W}_{u_1}, \quad (5a)$$

$$\frac{du_2}{dt} = (-d_u u_2 + \gamma v_f(t) u_2 + a u_1) + \sigma_u \dot{W}_{u_2}. \quad (5b)$$

Albeit linear, the model in (5) is still able to generate intermittent instability due 3  
to the time-periodic damping  $v_f(t) = f_0 + f_t \sin(\omega_f t + \phi)$  because the coefficient 4  
 $-d_u + \gamma v_f(t)$  switches between negative (stable) and positive (unstable) intervals 5  
during a time period. Therefore, as shown in panel (c) of Figure 7, the time- 6  
averaged PDF associated with the linear model (5) succeeds in capturing the 7  
non-Gaussian fat tailed PDF of the observed data given the optimal parameters 8  
in Table 1. However, obvious discrepancies in the autocorrelation functions 9  
and the power spectrums are observed in panel (a), (b) and (d) of Figure 7. 10  
Although the oscillations in the autocorrelation functions of the linear model 11  
are of the same frequency as that of the observed data, the decaying rate of the 12  
autocorrelation functions of the linear model is linear and much slower than the 13  
nonlinear decay rate of those associated with the BSISO indices, indicating a 14  
barrier in describing the nature of the BSISO indices by the linear model. In 15  
addition, the power spectrums of the linear model with only the additive noise 16  
are more concentrated compared with that of the observed data, implies the 17  
insufficient with only additive noise. Therefore, nonlinearity and multiplicative 18  
noise are necessary in order to match the statistics of the BSISO indices. 19

Column (c) in Figure 8 shows the skill scores for ensemble mean prediction 20  
utilizing the linear model (5) with the optimal parameters. The linear model has 21  
a comparably high prediction skill as the nonlinear low-order stochastic model 22  
(1) in the short and medium range up to 20-25 days (column (a) in Figure 23  
8). Above this range, the linear model becomes less skillful than the nonlinear 24  
model in years 1998, 2003, 2004 and 2005. On the other hand, the linear model 25  
shows the same skill with respect to the RMS error as the nonlinear model and 26  
even slightly more skillful prediction with respect to the bivariate correlation at 27  
lead times of 50-60 days in year 2001 and 2002 with regular oscillations. This 28  
possibly comes from the linear nature of the intraseasonal oscillation in these 29  
two years. 30

1 To further understand the difference in prediction utilizing the linear and  
 2 nonlinear models with optimal parameters, the predictions in year 2005 are  
 3 shown in column (a) and (c) of Figure 9. The forecasts at lead times of 25 and  
 4 35 days are shown in row I and II. The prediction at a lead time of 25 days  
 5 in June and July utilizing the nonlinear model has almost the same pattern as  
 6 the truth while that utilizing the linear model is slightly shifted forward in time  
 7 and therefore results in a lower correlation. The superiority of the nonlinear  
 8 model in this medium range forecasting is due to the phase correction from  
 9 the stochastic part  $\omega_u$  shown in the fourth row of Figure 13 in Appendix C.  
 10 Similarly, the predicted amplitude around June 1 at lead times of both 25 and  
 11 35 days utilizing the nonlinear model being more close to that utilizing the  
 12 linear model attributes to the initial correction of the overall damping by the  
 13 stochastic part  $v$ , as is shown in the third row of Figure 13. The prediction  
 14 including the ensemble members starting from June 1, the beginning of the  
 15 active mature phase, is shown in row IV. Except slightly shifted forward in  
 16 time in the prediction utilizing the linear model, the two models have the  
 17 comparable good ensemble mean prediction skill up to 20 days. However, the  
 18 medium- and long-range forecast up to 5 months utilizing the linear models  
 19 shows not only a large bias in the ensemble mean forecasting but a small  
 20 ensemble spread as well, which fails to include the truth and thus indicates  
 21 a false uncertainty quantification in prediction. The misleading forecast of the  
 22 linear model is associated with the long memory (slow decaying autocorrelation)  
 23 and the single frequency dominated nature (concentrated power spectrum) of  
 24 the system. Actually, the opposite patterns in the truth and prediction of the  
 25 linear model during the period from the middle of July to the end of October  
 26 are good indicators of the nonlinear nature of the oscillations in year 2005.

27 Another important thing to check is the model sensitivity. The model error  
 28 dependence on the parameters is shown in Figure 12 in Appendix B. Clearly,  
 29 the nonlinear low-order stochastic model (1) is more robust with respect to  
 30 the parameter variations around the optimal values compared with the linear  
 31 stochastic model (5). Note that when  $\gamma v_f(t) > d_u$ , a slight increase in  $\gamma v_f(t)$   
 32 leads to an exponential increase in  $u_1$  and  $u_2$  in the linear model (5) and therefore  
 33 the linear model is sensitive to the parameter variations. On the other hand,  
 34 the energy-conserving nonlinear interaction plays a significant role in increasing  
 35 the robustness in the nonlinear low-order stochastic model (1). Large values in  
 36  $u_1$  and  $u_2$  in the nonlinear low-order stochastic model (1) strongly reduce the  
 37 stochastic damping  $v$  due to the nonlinear feedback  $-\gamma(u_1^2 + u_2^2)$  in (1c), which  
 38 then prevents the unbounded exponential increase in  $u_1$  and  $u_2$  and guarantees  
 39 the robustness of the model.

To check the sensitivity of the prediction skill with respect to the parameter 1  
 variations. We pick up the following random suboptimal parameters in each 2  
 ensemble for both the nonlinear and linear models, 3

$$\begin{aligned}
 \sigma_u &= \sigma_u^* + \mathcal{U}(0, 0.2), & d_u &= d_u^* + \mathcal{U}(-0.4, 0), \\
 \sigma_v &= \sigma_v^* + \mathcal{U}(-0.2, 0.2), & d_v &= d_v^* + \mathcal{U}(-0.2, 0.2), \\
 \sigma_\omega &= \sigma_\omega^* + \mathcal{U}(-0.2, 0.2), & d_\omega &= d_\omega^* + \mathcal{U}(-0.2, 0.2), \\
 f_t &= f_t^* + \mathcal{U}(0, 1), & \gamma &= \gamma^* + \mathcal{U}(0, 0.3),
 \end{aligned} \tag{6}$$

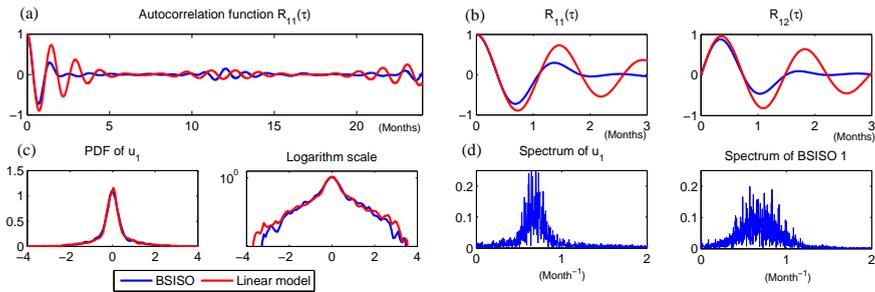
where the variables with asterisk are the optimal parameters and  $\mathcal{U}(a, b)$  is the 4  
 uniform distribution within the interval  $[a, b]$ . The forecasting skill scores are 5  
 shown in column (b) and (d) of Figure 8. As expected, due to the energy- 6  
 conserving nonlinear interactions, the nonlinear low-order stochastic model (1) 7  
 with the random suboptimal parameters has the comparable prediction skill 8  
 as that with the optimal parameters. On the other hand, the RMS error 9  
 in prediction utilizing the linear model (5) with the suboptimal parameters 10  
 increases dramatically as a function of lead time, contrasting the skillful 11  
 prediction with the optimal parameters. The overestimation in amplitude 12  
 utilizing the linear model (5) with the random suboptimal parameters at lead 13  
 times of 25 and 35 days shown in Figure 9 is a good evidence. In addition, with 14  
 the random suboptimal parameters, the ensemble spread for predicting starting 15  
 from April 1 overestimates the uncertainty at long ranges and that starting 16  
 from June 1 fails to include the truth, which together with the huge bias in the 17  
 ensemble mean estimation indicates the useless prediction. 18

As a remark, the nonlinear model (1) without the nonlinear feedback term 19  
 $-\gamma(u_1^2 + u_2^2)$  is the stochastic parameterized extended Kalman filter (SPEKF) 20  
 model [42, 43], which is able to capture the intermittent nature of the BSISO 21  
 indices and is skillful as a short-term forecast model. Yet, without the energy- 22  
 conserving nonlinear interactions, the SPEKF model is not robust for long-range 23  
 predictions and has sensitive dependence on parameters. 24

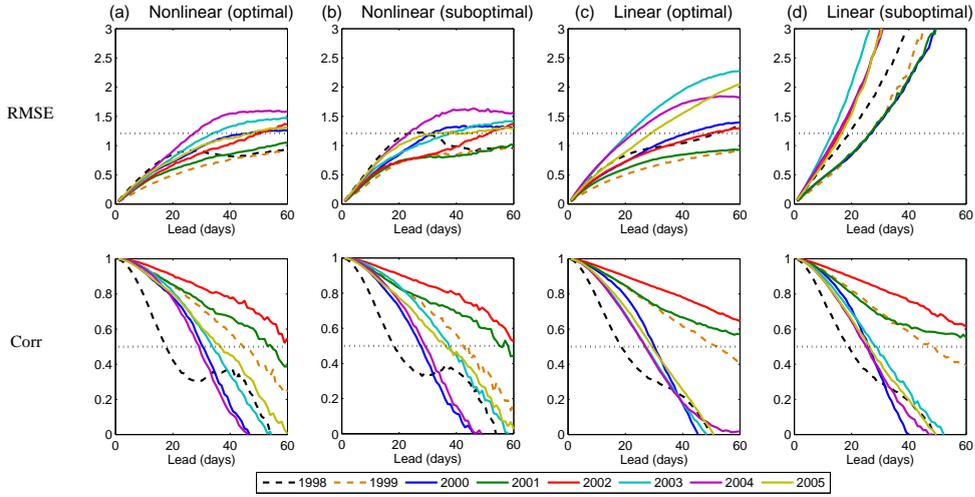
In summary, the linear model (5) with time-periodic damping  $v_f(t)$  is able 25  
 to capture the non-Gaussian fat tailed PDF of the truth and has comparably 26  
 high prediction skill as the nonlinear low-order stochastic model (1) in the years 27  
 with regular oscillations given the optimal parameters. The failure of the linear 28  
 model in capturing the autocorrelation functions and power spectrums leads to 29  
 a significant bias in medium- and long-range prediction with respect to ensemble 30  
 mean in the years with nonlinear oscillations and the ensemble spread fails 31  
 to include the truth in those years. Without the energy-conserving nonlinear 32  
 interaction, the linear model is also sensitive to the parameter variations around 33

1 the optimal values. The error in prediction increases dramatically utilizing the  
 2 linear model with the random suboptimal parameters.

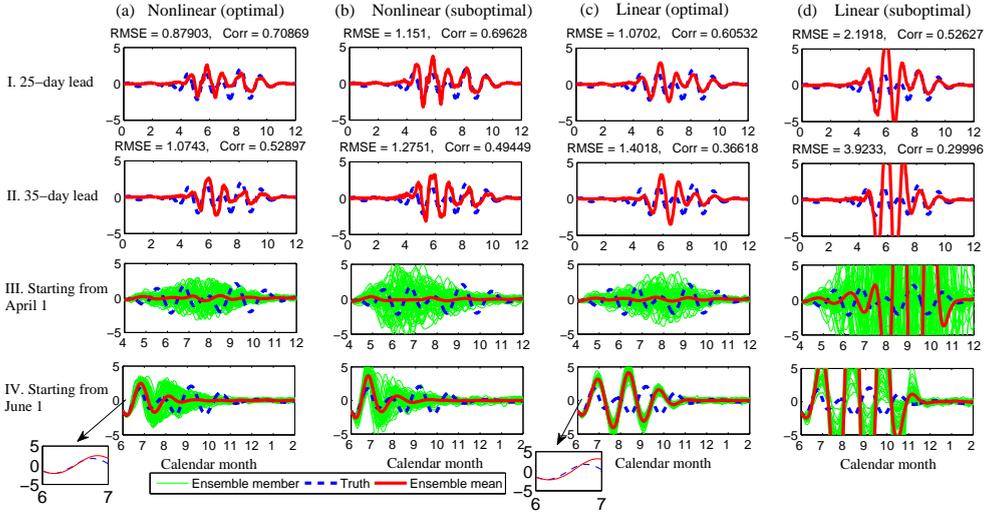
3 Note that the bivariate correlation in prediction utilizing the linear model  
 4 with the random suboptimal parameters remains almost the same as that  
 5 with the optimal parameters due to the fact that the linear oscillator model  
 6 is nevertheless able to capture the averaged oscillation frequency. Actually,  
 7 this is a representative example to demonstrate that the bivariate correlation  
 8 should not be overemphasized in measuring the prediction skill. In [33, 44],  
 9 an information-theoretic framework is proposed to evaluate the forecast skill,  
 10 which contains the information surrogates of the RMS error and the bivariate  
 11 correlation as well as the information deficiency in comparing the amplitudes  
 12 of the prediction and the truth. As shown in [33], two different predictions  
 13 can have nearly the same bivariate correlation with the truth but have quite  
 14 different information deficiency. Note that although only the RMS error and  
 15 bivariate correlation are utilized here in measuring the prediction skill of the  
 16 BSISO indices, the short- and medium-range ensemble mean prediction and  
 17 the long-term ensemble spread of the nonlinear low-order stochastic model (1)  
 18 have nearly the same amplitude as the truth, implying the small information  
 19 deficiency in prediction.



**Fig. 7.** Statistics of the linear stochastic model (5) with the optimal parameters in Table 1. (a) Long-term autocorrelation function  $R_{11}(\tau)$  from 0 to 24 months. (b) Short-term autocorrelation function  $R_{11}(\tau)$  and cross-correlation function  $R_{12}(\tau)$  from 0 to 3 months. (c) Equilibrium PDFs of the signal  $u_1$  from stochastic model compared with that of the BSISO 1 index. (d) Spectrum of  $u_1$  compared with that of the BSISO 1 index.



**Fig. 8.** Comparison of the skill scores with RMS error and bivariate correlation for prediction utilizing different models and parameters as a function of lead days in different years. Column (a): Nonlinear stochastic model (1) with optimal parameters; Column (b): Nonlinear stochastic model (1) with suboptimal parameters; Column (c): Linear stochastic model (5) with optimal parameters; Column (d): Linear stochastic model (1) with suboptimal parameters.



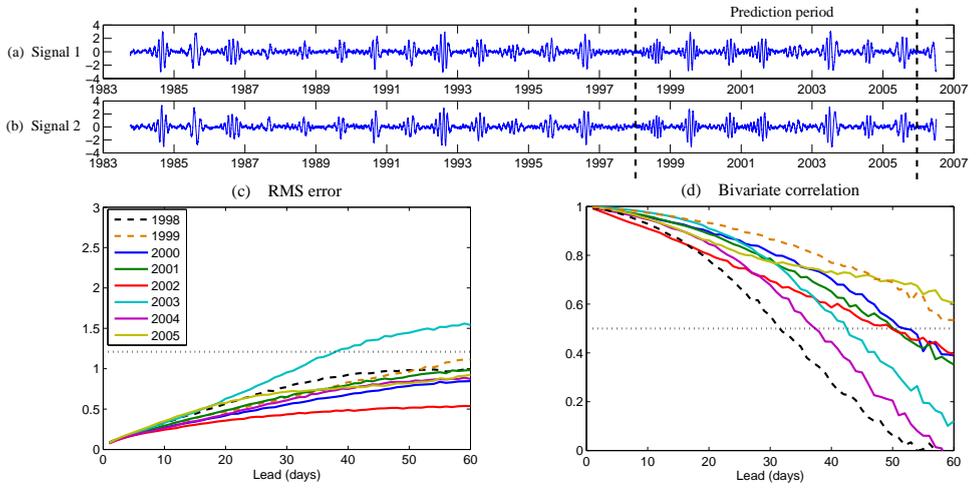
**Fig. 9.** Comparison of the prediction skill in year 2005 utilizing different models and parameters at lead times of 25 (Row I) and 35 (Row II) days and the long-term prediction up to 8 months starting from April 1 (Row III) and June 1 (Row IV). The skill scores are computed from BSISO 1 and 2 indices but only the curves of BSISO 2 index are shown. Column (a): Nonlinear stochastic model (1) with optimal parameters; Column (b): Nonlinear stochastic model (1) with suboptimal parameters; Column (c): Linear stochastic model (5) with optimal parameters; Column (d): Linear stochastic model (1) with suboptimal parameters.

## 6 Twin Experiment

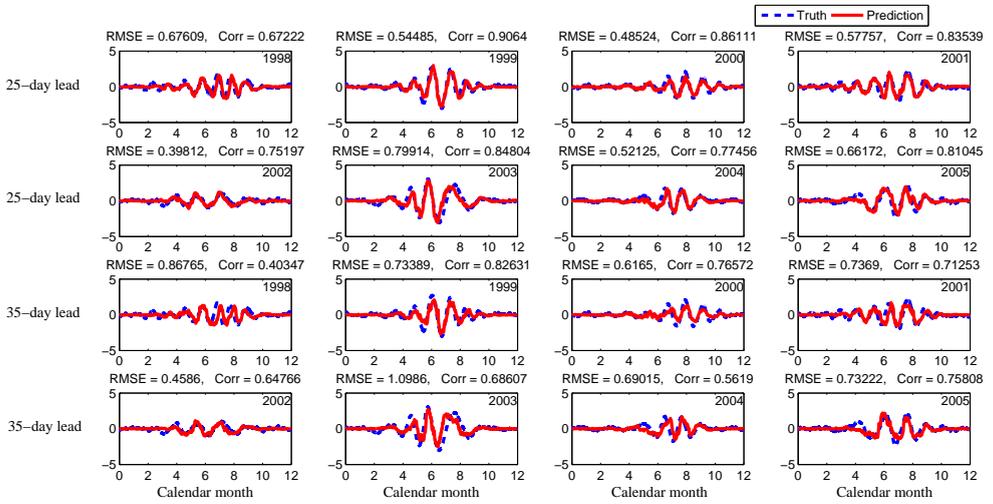
1

To explore the model error and the predictability limits of the nonlinear low- 2  
order stochastic model (1) in predicting the BSISO indices, we include the results 3  
from the perfect model twin experiment. In the twin experiment, the truth signal 4  
is generated from the nonlinear low-order stochastic model (1) and therefore the 5  
data assimilation algorithm in Section 3.2 is based on a perfect filter. The length 6  
of the two signals representing  $u_1$  and  $u_2$  in the twin experiment is the same as 7  
the two BSISO indices in Figure 1 and as in predicting the BSISO indices we 8  
predict the signals in the last eight years of these two time series. See panel (a) 9  
and (b) in Figure 10. 10

The prediction skill scores as a function of lead time are shown in panel 11  
(c) and (d) in Figure 10 and the ensemble mean predictions at lead times of 25 12  
and 35 days are shown in Figure 11. Actually, the internal prediction skill in 13  
this twin experiment indicates the predictability limit of the model. Comparing 14  
panel (c) and (d) in Figure 10 and panel (c) and (d) in Figure 3, it is clear 15  
that the BSISO prediction skill is comparable to this internal prediction skill, 16  
suggesting that the nonlinear low-order stochastic model (1) has a significant 17  
skill for determining the predicability limits of the large scale cloud patterns of 18  
the BSISO. 19



**Fig. 10.** Twin experiment. (Panel (a) and (b)) The two components of the signal generated from the nonlinear low-order stochastic model (1), where the eight year period from 1998 2006 is utilized as the prediction period. (Panel (c) and (d)) Skill scores with RMS error and bivariate correlation for prediction as a function of lead days in different years.



**Fig. 11.** Twin experiment. Prediction of the twin experiment signal 2 (panel (b) in Figure 10) generated from the nonlinear low-order stochastic model (1) regarding the ensemble mean at lead times of 25 and 35 days in different years. The blue curve shows the truth in panel (b) of Figure 10 and the red curve shows the ensemble average of the predicted signal. The number of ensemble utilized for prediction is 50.

## 1 7 Conclusions and Discussions

2 A recent developed nonlinear data analysis technique NLSA has been applied  
3 to the CLAUS data to define two spatial patterns associated with the BSISO  
4 without detrending or spatial-temporal filtering. This provides two BSISO  
5 indices with strong intermittency and non-Gaussian fat-tailed PDF (Figure  
6 1). Then a recent systematic strategy for data driven physics-constrained low-  
7 order stochastic modeling is applied to the BSISO indices. The result is a four  
8 dimensional nonlinear stochastic model (1) with two state variables denoting  
9 the observed BSISO indices and two hidden variables representing stochastic  
10 damping and stochastic phase. The model contains correlated multiplicative  
11 noise through the energy-conserving nonlinear interactions between the observed  
12 and hidden variables as well as the additive noise.

13 The parameters in the nonlinear low-order stochastic model (1) are calibrated  
14 in minimizing the model error of the time-averaged PDF compared with that  
15 of the truth in the information framework. The model with the optimal  
16 parameters succeeds in recovery the highly intermittent non-Gaussian PDF; the  
17 autocorrelations and the power spectrums of the model signals almost perfectly  
18 match those of the observed BSISO indices (Figure 2). The special structure  
19 of model (1) allows an effective data assimilation algorithm for determining the  
20 initial ensemble of the hidden variables in the ensemble forecasting scheme. This  
21 on-line prediction algorithm shows that the ensemble mean forecasting skill of  
22 the BSISO in the non-El Niño year is at least 30 days and even reaches 55 days  
23 in the years with regular intraseasonal oscillation. In the strong El Niño year  
24 (year 1998) the useful prediction is around 18 days due to both the increase of  
25 noise to signal ratio and the intrinsic model error (Figure 3–5). Furthermore,  
26 the ensemble spread is a good indicator of the forecasting uncertainty at long  
27 range (Figure 6).

28 The twin experiment (Figure 10 and 11) shows the skill of forecasts in the  
29 perfect modeling setting is comparable with that of predicting the BSISO indices,  
30 implying that the nonlinear low-order stochastic model (1) has a significant skill  
31 for determining the predicability limits of the large scale cloud patterns of the  
32 BSISO.

33 To check the role of the nonlinearity in the low-order models, the prediction  
34 skill of the linear model (5) is studied. With time-periodic damping  $v_f(t)$ , the  
35 linear model is able to capture the non-Gaussian fat tailed PDF of the truth  
36 and has comparably high prediction skill as the nonlinear low-order stochastic  
37 model (1) in the years with regular oscillations given the optimal parameters.  
38 The failure of the linear model in capturing the autocorrelation functions and

power spectrums (Figure 7) leads to a significant bias in medium- and long-range 1 prediction with respect to ensemble mean in the years with nonlinear oscillations 2 and the ensemble spread fails to include the truth in those years (Figure 8 and 3 9). In addition, the linear model is sensitive to the parameter variations around 4 the optimal values (Figure 12). The error in prediction increases dramatically 5 utilizing the linear model with random suboptimal parameters (Figure 8 and 9). 6 On the other hand, the SPEKF model, which is a nonlinear model but has no 7 energy-conserving nonlinear interaction, is also quite sensitive with respect to 8 the parameter variations for long-range forecasting. 9

Although in this paper, we focus only on the skill of predicting the two 10 BSISO time series, it is straightforward to translate to the prediction of the 11 location and evolution of BSISO convection itself. Recall that the original 12 spatiotemporal patterns of the MJO cloud clusters are illustrated in the video 13 in [24]. The predicted spatial patterns are a rank-2 reconstruction constructed 14 from the predicted temporal patterns ( $u_1$  and  $u_2$ ) multiplied by the dataset 15 projected onto the original temporal patterns. Because  $u_1$  and  $u_2$  evolve in 16 near-quadrature (as do the original temporal patterns), we do not expect major 17 qualitative differences between the structure of the predicted cloud clusters in 18 the reconstruction and the original clusters when the time series can be predicted 19 reasonably well. 20

Most of the current studies of BSISO prediction is up to 30 days. In [45], a 21 regression scheme is designed to study the forecasts of central India precipitation. 22 The prediction skill lasts for 30 days with respect to the pattern correlation. 23 Yet, their prediction underestimates the amplitudes and fails to predict some 24 monsoon onsets. Similar results are found in [46] that the skill utilizing the 25 coupled ocean-atmosphere forecast models for monsoon prediction is up to 30 26 days as well. Thus, our nonlinear low-order stochastic model combined with 27 NLSA data analysis tool are potentially able to extend the prediction limit of 28 BSISO and boreal summer monsoon. 29

Note that what we have predicted utilizing the nonlinear low-order stochastic 30 model (1) is the anomalies of the BSISO. Actually, when we apply NLSA to the 31 CLAU data, besides the BSISO modes, we also obtain those large-scale modes 32 such as the annual and seasonal modes, which reflect the background of rainfall. 33 These modes are quite regular and easy to predict. Therefore, to recover the 34 total rainfall, we simply need to combine the predicted BSISO anomaly with 35 the background. 36

We have also looked at the prediction skill utilizing the nonlinear low-order 37 stochastic model (1) in predicting the boreal winter MJO derived from NLSA 38 with shift map and with the improved prediction scheme incorporating the 39 effective data assimilation algorithm in Section 3.2. The boreal winter MJO time 40

1 series also have non-Gaussian fat tailed time-averaged PDF and the oscillation  
2 frequency is slightly lower than that of the BSISO. Yet, unlike BSISO with  
3 moderate or strong amplitudes in almost all the years, the amplitudes of the  
4 boreal winter MJO range from weak to strong in different years. The forecasting  
5 skill for the BSISO and the boreal winter MJO is comparable in the moderate  
6 and strong years while the small signal to noise ratio in the weak MJO years,  
7 including year 1998, deteriorates the prediction skill. The results are similar to  
8 those reported in [32]

9 We point out that the information-theoretic calibration procedure utilized  
10 in this work is widely adopted as training strategy for improving the predictive  
11 skill in many different issues. Imperfect predictions via Multi Model Ensemble  
12 forecasts are improved with the information-theoretic framework [47]. The  
13 prediction skill of imperfect large-dimensional turbulent models are enhanced  
14 through statistical response and information theory [48]. The forecasting  
15 skill of the RMM index is also greatly enhanced by combining three different  
16 information measures compared with adopting only path-wise measures [33, 44].

17 We also note that the conditional Gaussian models as in (11) and the  
18 effective data assimilation algorithm in Section 3.2 are innovative tools for  
19 studying the hidden processes from the observations in the turbulent flows. One  
20 example is the understanding of the practical information barrier in recovering  
21 the fluid flows with noisy Lagrangian tracers [49, 50], which contains the  
22 application of the information theory as well.

23 **Acknowledgment:** This research of A.J.M is partially supported by the Office  
24 of Naval Research grant ONR MURI N00014-12-1-0912. N.C. is supported as a  
25 graduate research assistant on this grant. A.J.M also gratefully acknowledges the  
26 financial support given by the Earch System Science Organization, Ministry of  
27 Earth Sciences, Government of India (Grant no./Project no MM/SERP/CNRS/  
28 2013/INT-10/002) to conduct this research under Monsoon Mission. The  
29 authors thank Dimitrios Giannakis and Eniko Szekely for discussion on NLSA  
30 applied to CLAUS data.

## 8 Appendix 1

### 8.1 A. Reconstruction of the eigenfunctions from the shift map in NLSA 2 3

Following Section 2 in [24], the Laplace-Beltrami operator  $L$  is given by formula (3) there. In formula (4) there, the eigenfunction  $\phi_i$  with  $i \in \{0, 1, 2, \dots\}$  is solved via 4  
5  
6

$$L\phi_i = \lambda_i\phi_i, \quad (7)$$

where  $\lambda_i$  is the associated eigenvalue. Now we utilize the shift map  $S_\tau$  [35], where  $S_\tau f(x(t_i)) = f(x(t_{i+\tau}))$ , to formulate a new operator  $A$ ,

$$A_{ij} = \langle \phi_i, S_\tau \phi_j \rangle.$$

Then we solve the following eigenvalue problem 7

$$Av_i = \tilde{\lambda}_i v_i, \quad (8)$$

where the eigenvectors  $v_i$  are utilized as the time series of different modes. Similar to the formula (5) in [24] but replacing  $\phi_i$  by  $v_i$  that the spatiotemporal patterns are recovered through the operation 8  
9  
10

$$\tilde{X}_i = XDv_i v_i^T, \quad (9)$$

where  $X$  is the lagged embedding of the raw data and  $D$  is the weight matrix. 11

To derive the BSISO indices in Figure 1,  $\tau$  is set to be 1. 12

### 8.2 B. Calibration of the nonlinear and linear low-order stochastic models and sensitivity analysis 13 14

The optimal parameters in the nonlinear low-order stochastic model (1) and the linear model (5) are calibrated by systematically minimizing the information distance, i.e., the model error, of the time-averaged PDF of the model  $\pi^M$  compared with that of the BSISO index  $\pi$  [38, 39, 51, 52, 53], 15  
16  
17  
18

$$\mathcal{P}(\pi, \pi^M) = \int \pi \log \left( \frac{\pi}{\pi^M} \right). \quad (10)$$

The model error dependence on the variation of different parameters is shown in Figure 12, which indicates that the nonlinear low-order stochastic model (1) 19  
20

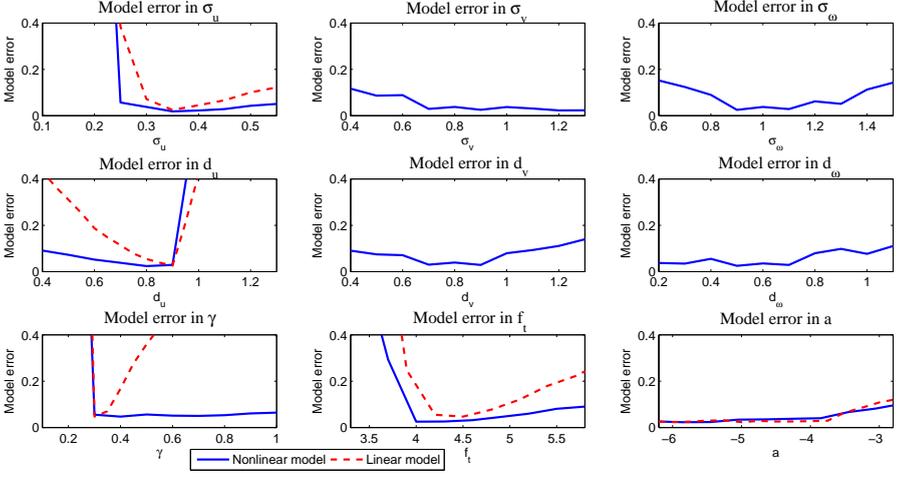
1 is robust with respect to the parameters around their optimal values. The huge  
 2 model error with the underestimation of  $\sigma_u, f_t$  and  $\gamma$  and the overestimation of  
 3  $d_u$  is due to the failure of capturing the intermittency. Note that the model error  
 4 has only a weak dependence on the background phase  $a$  since the contribution  
 5 of the oscillation in the signal has been averaged out in the time-averaged PDF.  
 6 However, the parameter  $a$  is crucial in describing the frequency of intraseasonal  
 7 oscillation, and it is calibrated by matching the autocorrelation functions  
 8 associated with the model and the truth. The other parameters  $d_v, \sigma_v, d_\omega$  and  
 9  $\sigma_\omega$  in describing the stochastic processes affect not only on the model error but  
 10 more on the autocorrelations and power spectrums as well. A large discrepancy  
 11 appears in the statistics if these parameters are outside the optimal range. The  
 12 parameter  $f_0$  is not an independent parameter given  $d_u$  and  $\gamma$  and therefore we  
 13 fix its value. The frequency  $\omega_f$  in the time-periodic damping  $v_f(t)$  is prescribed  
 14 to be  $2\pi/12$  such that one time unit of the model corresponds to one month in  
 15 reality. The phase  $\phi$  in  $v_f(t)$  is tuned to make the strong intermittency occur  
 16 in the boreal summer in accordance with the BSISO indices. Note that none  
 17 of the parameters is redundant in the nonlinear stochastic model (1). In fact,  
 18 without the hidden variables  $v$  and  $\omega_u$ , even if the time-period damping  $v_f(t)$  is  
 19 able to crudely describe the active phase of BSISO in the reduced linear model,  
 20 a distinguished disparity is observed in the model statistics compared with the  
 21 truth, indicating the intrinsic barrier [39, 52]. In addition, as seen in Figure 12,  
 22 the linear model (5) is more sensitive with respect to the parameter variations  
 23 around the optimal values.

### 24 **8.3 C. Mathematical details of effective data assimilation** 25 **and prediction algorithm**

Recall the nonlinear low-order stochastic model (1). Denote by  $\mathbf{U} = (u_1, u_2)^T$   
 and  $\mathbf{\Gamma} = (v, \omega_u)^T$ . The abstract form of the low-order stochastic model (1) is  
 given as follows:

$$d\mathbf{U}_t = [\mathbf{A}_0(t, \mathbf{U}) + \mathbf{A}_1(t, \mathbf{U})\mathbf{\Gamma}_t]dt + \mathbf{\Sigma}_U(t, \mathbf{U})d\mathbf{W}_U(t), \quad (11a)$$

$$d\mathbf{\Gamma}_t = [\mathbf{a}_0(t, \mathbf{U}) + \mathbf{a}_1(t, \mathbf{U})\mathbf{\Gamma}_t]dt + \mathbf{\Sigma}_\Gamma(t, \mathbf{U})d\mathbf{W}_\Gamma(t), \quad (11b)$$



**Fig. 12.** Sensitivity analysis of the parameters in the nonlinear low-order stochastic model(1) (solid blue curves) and the linear stochastic model (5) (dashed red curves).

where

$$\mathbf{A}_0 = \begin{pmatrix} -d_u u_1 + \gamma v_f(t) u_1 - a u_2 \\ -d_u u_2 + \gamma v_f(t) u_2 + a u_1 \end{pmatrix}, \quad \mathbf{A}_1 = \begin{pmatrix} \gamma u_1 & -u_2 \\ \gamma u_2 & u_1 \end{pmatrix},$$

$$\mathbf{a}_0 = \begin{pmatrix} -\gamma(u_1^2 + u_2^2) \\ 0 \end{pmatrix}, \quad \mathbf{a}_1 = \begin{pmatrix} -d_v \\ -d_\omega \end{pmatrix},$$

$$\Sigma_U = \begin{pmatrix} \sigma_u & \\ & \sigma_u \end{pmatrix}, \quad \Sigma_\Gamma = \begin{pmatrix} \sigma_v & \\ & \sigma_\omega \end{pmatrix}.$$

The model (11) is a conditional Gaussian system conditioned on the observations 1  
 $\mathbf{U}$ , meaning that once the observations  $\mathbf{U}$  are given the dynamics of  $\Gamma$  in (11) 2  
becomes a Gaussian system. The special structure of system (11) allows 3  
the closed analytic formulae for the evolution of the conditional Gaussian 4  
distributions of the hidden parameters  $v$  and  $\omega_u$  [40] obtained in the Bayesian 5  
framework: 6

$$d\mu_t = [\mathbf{a}_0(t, \mathbf{U}) + \mathbf{a}_1(t, \mathbf{U})\mu_t]dt + (R_t \mathbf{A}_1^*(t, \mathbf{U}))(\Sigma_U \Sigma_U^*)^{-1}(t, \mathbf{U}) \times$$

$$[d\mathbf{U}_t - (\mathbf{A}_0(t, \mathbf{U}) + \mathbf{A}_1(t, \mathbf{U})\mu_t)dt],$$

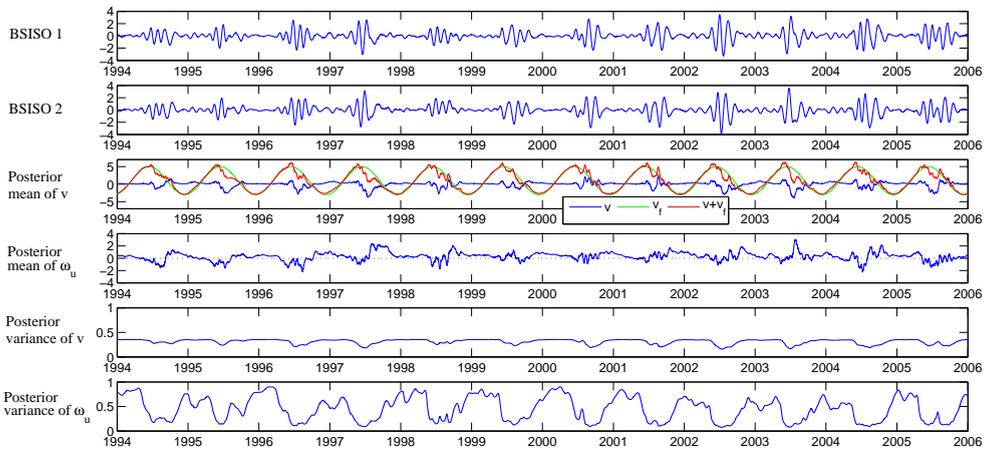
$$dR_t = \{\mathbf{a}_1(t, \mathbf{U})R_t + R_t \mathbf{a}_1^*(t, \mathbf{U}) + (\Sigma_\Gamma \Sigma_\Gamma^*)(t, \mathbf{U})$$

$$-(R_t \mathbf{A}_1^*(t, \mathbf{U}))(\Sigma_U \Sigma_U^*)^{-1}(t, \mathbf{U})(R_t \mathbf{A}_1^*(t, \mathbf{U}))^*\} dt, \quad (12)$$

1 where  $\mu_t$  and  $R_t$  are the posterior mean and posterior covariance of the  
 2 conditional distributions, respectively. The asterisk represents the complex  
 3 conjugate.

4 As a remark, the formulae (12) are optimal if and only if the signal is  
 5 generated from system (11). Since our observed signal, i.e., the BSISO indices,  
 6 are not from the nonlinear low-order stochastic model (11), the evolutions of the  
 7 conditional Gaussian distributions (12) are suboptimal.

8 In Figure 13 we show the posterior mean and variance of stochastic damping  
 9  $v$  and stochastic phase  $\omega_u$  in (1) as a function of time compared with the  
 10 observations of the BSISO indices from January 1994 to December 2005.  
 11 Note that the corrections of the stochastic variables  $v$  and  $\omega_u$  from the data  
 12 assimilation algorithm are significant and important in the intermittent phase.  
 13 In addition, the posterior covariance at the intermittent phase is smaller than  
 14 that at the quiescent phase, indicating the small uncertainty in the recovered  
 15 stochastic variables.



**Fig. 13.** Recovery of posterior mean and variance of stochastic damping  $v$  and stochastic phase  $\omega_u$  in (1) as a function of time compared with the observations of the BSISO indices from January 1994 to December 2005. The cross-covariance of  $v$  and  $\omega_u$  is negligible of order  $O(10^{-18})$  and is not shown here.

## References

- 1
- [1] William KM Lau and Duane E Waliser. *Intraseasonal variability in the atmosphere-ocean climate system*. Springer, 2012. 2
- [2] Peter J Webster, Vo Oo Magana, TN Palmer, J Shukla, RA Tomas, M u Yanai, and T Yasunari. Monsoons: Processes, predictability, and the prospects for prediction. *Journal of Geophysical Research: Oceans (1978–2012)*, 103(C7):14451–14510, 1998. 3  
4  
5  
6
- [3] Tiruvalam Natarajan Krishnamurti and D Subrahmanyam. The 30-50 day mode at 850 mb during MONEX. *Journal of the Atmospheric Sciences*, 39(9):2088–2095, 1982. 7  
8
- [4] In-Sik Kang, Chang-Hoi Ho, Young-Kwon Lim, and KM Lau. Principal modes of climatological seasonal and intraseasonal variations of the Asian summer monsoon. *Monthly weather review*, 127(3):322–340, 1999. 9  
10  
11
- [5] Qinghua Ding and Bin Wang. Predicting extreme phases of the Indian summer monsoon. *Journal of Climate*, 22(2):346–363, 2009. 12  
13
- [6] V Krishnamurthy and J Shukla. Intraseasonal and seasonally persisting patterns of Indian monsoon rainfall. *Journal of climate*, 20(1):3–20, 2007. 14  
15
- [7] V Krishnamurthy and J Shukla. Seasonal persistence and propagation of intraseasonal patterns over the Indian monsoon region. *Climate Dynamics*, 30(4):353–369, 2008. 16  
17
- [8] Bin Wang, June-Yi Lee, In-Sik Kang, J Shukla, C-K Park, A Kumar, J Schemm, S Cocks, J-S Kug, J-J Luo, et al. Advance and prospectus of seasonal prediction: assessment of the APCC/CLIPAS 14-model ensemble retrospective seasonal prediction (1980–2004). *Climate Dynamics*, 33(1):93–117, 2009. 18  
19  
20  
21
- [9] June-Yi Lee, Bin Wang, I-S Kang, J Shukla, A Kumar, J-S Kug, JKE Schemm, J-J Luo, T Yamagata, X Fu, et al. How are seasonal prediction skills related to models' performance on mean state and annual cycle? *Climate dynamics*, 35(2-3):267–283, 2010. 22  
23  
24  
25
- [10] In-Sik Kang, June-Yi Lee, and Chung-Kyu Park. Potential predictability of summer mean precipitation in a dynamical seasonal prediction system with systematic error correction. *Journal of climate*, 17(4):834–844, 2004. 26  
27  
28
- [11] Bin Wang, Qinghua Ding, Xiouhua Fu, In-Sik Kang, Kyung Jin, J Shukla, and Francisco Doblas-Reyes. Fundamental challenge in simulation and prediction of summer monsoon rainfall. *Geophysical Research Letters*, 32(15), 2005. 29  
30  
31
- [12] Hye-Mi Kim and In-Sik Kang. The impact of ocean–atmosphere coupling on the predictability of boreal summer intraseasonal oscillation. *Climate dynamics*, 31(7-8):859–870, 2008. 32  
33  
34
- [13] DR Pattanaik and Arun Kumar. Prediction of summer monsoon rainfall over India using the NCEP climate forecast system. *Climate dynamics*, 34(4):557–572, 2010. 35  
36
- [14] Nachiketa Acharya, Sarat C Kar, UC Mohanty, Makarand A Kulkarni, and SK Dash. Performance of GCMs for seasonal prediction over india – a case study for 2009 monsoon. *Theoretical and applied climatology*, 105(3-4):505–520, 2011. 37  
38  
39
- [15] Archana Nair, UC Mohanty, Andrew W Robertson, TC Panda, Jing-Jia Luo, and Toshio Yamagata. An analytical study of hindcasts from general circulation models for Indian summer monsoon rainfall. *Meteorological Applications*, 21(3):695–707, 2014. 40  
41  
42  
43

- 1 [16] Sun-Seon Lee, Bin Wang, Duane E Waliser, Joseph Mani Neena, and June-Yi Lee.  
2 Predictability and prediction skill of the boreal summer intraseasonal oscillation in the  
3 intraseasonal variability hindcast experiment. *Climate Dynamics*, pages 1–13, 2015.
- 4 [17] Timothy DelSole and J Shukla. Linear prediction of indian monsoon rainfall. *Journal*  
5 *of Climate*, 15(24):3645–3658, 2002.
- 6 [18] Charles Jones, Leila MV Carvalho, R Wayne Higgins, Duane E Waliser, and JK E  
7 Schemm. A statistical forecast model of tropical intraseasonal convective anomalies.  
8 *Journal of climate*, 17(11):2078–2095, 2004.
- 9 [19] M Rajeevan, DS Pai, R Anil Kumar, and B Lal. New statistical models for long-range  
10 forecasting of southwest monsoon rainfall over india. *Climate Dynamics*, 28(7-8):813–  
11 828, 2007.
- 12 [20] Sun-Seon Lee and Bin Wang. Regional boreal summer intraseasonal oscillation  
13 over indian ocean and western pacific: comparison and predictability study. *Climate*  
14 *Dynamics*, pages 1–17, 2015.
- 15 [21] Matthew C Wheeler and Harry H Hendon. An all-season real-time multivariate MJO  
16 index: Development of an index for monitoring and prediction. *Monthly Weather*  
17 *Review*, 132(8):1917–1932, 2004.
- 18 [22] June-Yi Lee, Bin Wang, Matthew C Wheeler, Xiouhua Fu, Duane E Waliser, and  
19 In-Sik Kang. Real-time multivariate indices for the boreal summer intraseasonal  
20 oscillation over the Asian summer monsoon region. *Climate Dynamics*, 40(1-2):493–  
21 509, 2013.
- 22 [23] E Suhas, JM Neena, and BN Goswami. An indian monsoon intraseasonal oscillations  
23 (miso) index for real time monitoring and forecast verification. *Climate dynamics*,  
24 40(11-12):2605–2616, 2013.
- 25 [24] Eniko Székely, Dimitrios Giannakis, and Andrew J Majda. Extraction and predictability  
26 of coherent intraseasonal signals in infrared brightness temperature data. *Climate Dyn*,  
27 2014.
- 28 [25] Dimitrios Giannakis, Wen-wen Tung, and Andrew J Majda. Hierarchical structure  
29 of the Madden-Julian oscillation in infrared brightness temperature revealed through  
30 nonlinear Laplacian spectral analysis. In *Intelligent Data Understanding (CIDU), 2012*  
31 *Conference on*, pages 55–62. IEEE, 2012.
- 32 [26] Dimitrios Giannakis and Andrew J Majda. Comparing low-frequency and intermittent  
33 variability in comprehensive climate models through nonlinear Laplacian spectral  
34 analysis. *Geophysical Research Letters*, 39(10), 2012. doi: 10.1029/2012GL051575.
- 35 [27] Dimitrios Giannakis and Andrew J Majda. Nonlinear Laplacian spectral analysis  
36 for time series with intermittency and low-frequency variability. *Proceedings of the*  
37 *National Academy of Sciences*, 109(7):2222–2227, 2012.
- 38 [28] Dimitrios Giannakis and Andrew J Majda. Nonlinear Laplacian spectral analysis:  
39 capturing intermittent and low-frequency spatiotemporal patterns in high-dimensional  
40 data. *Statistical Analysis and Data Mining*, 6(3):180–194, 2013.
- 41 [29] Wen-wen Tung, Dimitrios Giannakis, and Andrew J Majda. Symmetric and  
42 antisymmetric convection signals in the Madden-Julian oscillation. Part I: basic modes  
43 in infrared brightness temperature. *Journal of the Atmospheric Sciences*, 71(9):3302–  
44 3326, 2014.
- 45 [30] Andrew J Majda and John Harlim. Physics constrained nonlinear regression models for  
46 time series. *Nonlinearity*, 26(1):201–217, 2013.

- [31] John Harlim, Adam Mahdi, and Andrew J Majda. An ensemble Kalman filter for statistical estimation of physics constrained nonlinear regression models. *Journal of Computational Physics*, 257:782–812, 2014.
- [32] Nan Chen, Andrew J Majda, and Dimitrios Giannakis. Predicting the cloud patterns of the Madden-Julian oscillation through a low-order nonlinear stochastic model. *Geophysical Research Letters*, 41(15):5612–5619, 2014.
- [33] Nan Chen and Andrew J Majda. Predicting the real-time multivariate Madden-Julian oscillation index through a low-order nonlinear stochastic model. *Monthly Weather Review*, 2015. in press.
- [34] KI Hodges, DW Chappell, GJ Robinson, and G Yang. An improved algorithm for generating global window brightness temperatures from multiple satellite infrared imagery. *Journal of Atmospheric & Oceanic Technology*, 17(10):1296–1312, 2000.
- [35] Tyrus Berry, Dimitrios Giannakis, and John Harlim. Nonparametric forecasting of low-dimensional dynamical systems. *arXiv preprint arXiv:1411.5069*, 2014.
- [36] S Kravtsov, D Kondrashov, and M Ghil. Multilevel regression modeling of nonlinear processes: Derivation and applications to climatic variability. *Journal of Climate*, 18(21):4404–4424, 2005.
- [37] D Kondrashov, MD Chekroun, AW Robertson, and M Ghil. Low-order stochastic model and “past-noise forecasting” of the Madden-Julian Oscillation. *Geophysical Research Letters*, 40(19):5305–5310, 2013.
- [38] Andrew J Majda and Boris Gershgorin. Quantifying uncertainty in climate change science through empirical information theory. *Proceedings of the National Academy of Sciences*, 107(34):14958–14963, 2010.
- [39] Andrew J Majda and Boris Gershgorin. Improving model fidelity and sensitivity for complex systems through empirical information theory. *Proceedings of the National Academy of Sciences*, 108(25):10044–10049, 2011.
- [40] Robert S Liptser and Albert N Shiryaev. *Statistics of Random Processes II: II. Applications*, volume 2. Springer, 2001.
- [41] M Berkelhammer, A Sinha, M Mudelsee, H Cheng, K Yoshimura, and J Biswas. On the low-frequency component of the enso–indian monsoon relationship: a paired proxy perspective. *Climate of the Past*, 10(2):733–744, 2014.
- [42] Boris Gershgorin, John Harlim, and Andrew J Majda. Improving filtering and prediction of spatially extended turbulent systems with model errors through stochastic parameter estimation. *Journal of Computational Physics*, 229(1):32–57, 2010.
- [43] Boris Gershgorin, John Harlim, and Andrew J Majda. Test models for improving filtering with model errors through stochastic parameter estimation. *Journal of Computational Physics*, 229(1):1–31, 2010.
- [44] M Branicki and AJ Majda. Quantifying bayesian filter performance for turbulent dynamical systems through information theory. *Communications in Mathematical Sciences*, 12(5), 2014.
- [45] Peter J Webster and Carlos Hoyos. Prediction of monsoon rainfall and river discharge on 15–30-day time scales. *Bulletin of the American Meteorological Society*, 85(11):1745–1765, 2004.
- [46] MJ MCPHADEN. Rama: The research moored array for african-asian-australian monsoon analysis and prediction. *Bull. Amer. Meteor. Soc.*, 90:459–480, 2009.

- 1 [47] Michal Branicki and Andrew J Majda. An information-theoretic framework for  
2 improving imperfect predictions via Multi Model Ensemble forecasts. *J. Nonlinear*  
3 *Science*, 2014. in press.
- 4 [48] Andrew J Majda and Qi Di. Improving prediction skill of imperfect turbulent models  
5 through statistical response and information theory. *J. Nonlinear Science*, 2015.  
6 submitted.
- 7 [49] Nan Chen, Andrew J Majda, and Xin T Tong. Information barriers for noisy  
8 Lagrangian tracers in filtering random incompressible flows. *Nonlinearity*, 27(9):2133–  
9 2163, 2014.
- 10 [50] Nan Chen, Andrew J Majda, and Xin T Tong. Noisy Lagrangian tracers for filtering  
11 random rotating compressible flows. *Journal of Nonlinear Science*, 2015. in press.
- 12 [51] Richard Kleeman. Measuring dynamical prediction utility using relative entropy.  
13 *Journal of the atmospheric sciences*, 59(13):2057–2072, 2002.
- 14 [52] Andrew J Majda and Michal Branicki. Lessons in uncertainty quantification for  
15 turbulent dynamical systems. *Discrete Contin. Dyn. Syst.*, 32(9):3133–3231, 2012.
- 16 [53] Michal Branicki, Nan Chen, and Andrew J Majda. Non-Gaussian test models for  
17 prediction and state estimation with model errors. *Chinese Annals of Mathematics,*  
18 *Series B*, 34(1):29–64, 2013.