



<sup>17</sup> \**Corresponding author address:* Department of Mathematics and Center for Atmosphere Ocean  
<sup>18</sup> Science, Courant Institute of Mathematical Sciences, New York University, 251 Mercer Street,  
<sup>19</sup> New York, NY 10012, USA.  
<sup>20</sup> E-mail: [chennan@cims.nyu.edu](mailto:chennan@cims.nyu.edu)

## ABSTRACT

21 We assess the predictability of the Monsoon Intraseasonal Oscillations (MIS-  
22 Os) as measured by precipitation. A recent advanced nonlinear time series  
23 technique, Nonlinear Laplacian Spectral Analysis, is applied directly to the  
24 daily rainfall data without any preliminary detrending or spatiotemporal fil-  
25 tering to define two spatial modes associated with the MISO. The time series of  
26 these two modes are highly intermittent with large variation in amplitude from  
27 year to year in the boreal summer season. Then a recent systematic strategy  
28 for data driven physics constrained low-order stochastic modeling is applied to  
29 these time series. The result is a four dimensional nonlinear stochastic model  
30 for the two observed MISO variables as well as two hidden variables involv-  
31 ing correlated multiplicative noise defined through energy conserving nonlinear  
32 interaction. Systematic calibration and prediction experiments with the nonlin-  
33 ear stochastic model show that the precipitation MISO indices can be skillfully  
34 predicted 20 to 50 days in advance and the ensemble spread in the forecast  
35 model is an accurate indicator of forecast uncertainty at long lead times. Then  
36 an effective and practical spatiotemporal reconstruction algorithm is developed,  
37 which shows the predicted spatiotemporal patterns have comparable skill as the  
38 MISO indices. It is also found that a 3-year short training period is sufficient  
39 for the model to describe the essential characteristics of the MISO and retain  
40 skillful predictions. In addition, outgoing longwave radiation is shown to be a  
41 good proxy for monsoon intraseasonal precipitation and the lagged embedding  
42 window size is crucial to reaching unbiased MISO indices.

## 43 **1. Introduction**

44 Monsoon Intraseasonal Oscillation (MISO) (Kikuchi et al. 2012; Lee et al. 2013; Sikka and  
45 Gadgil 1980; Goswami and Mohan 2001; Lau and Waliser 2011; Webster et al. 1998) is one of  
46 the prominent modes of tropical intraseasonal variability. As a slow moving planetary scale en-  
47 velope of convection propagating northeastward, it is strongly associated with the boreal summer  
48 monsoon rainfall over south Asia. Due to the interaction with the mean monsoon circulation and  
49 other modes of tropical variability, the propagating characteristics of the MISO are more complex  
50 compared with the eastward-propagating Madden-Julian Oscillation (MJO) (Zhang 2005). The  
51 MISO plays an important role in determining the onset and demise of the Indian summer mon-  
52 soon as well as affecting the rainfall over the Indian subcontinent (Murakami et al. 1986; Goswami  
53 and Mohan 2001; Goswami et al. 2003; Gadgil 2003). Therefore, the extended range prediction  
54 of MISO phases and real-time monitoring of the MISO have large societal impacts (Sahai et al.  
55 2013; Abhilash et al. 2014a).

56 Several indices have been proposed for real-time monitoring and forecast of the MISO. The  
57 Indian Institute of Tropical Meteorology (IITM) relies on an index based on extended empirical  
58 orthogonal function (EEOF) analysis on longitudinal averaged daily rainfall anomalies for extend-  
59 ed range prediction of MISO (Suhas et al. 2013; Sahai et al. 2013; Abhilash et al. 2014a). Another  
60 well-known MISO index (Lee et al. 2013) mimics that for the real-time multivariate MJO (RM-  
61 M) index (Wheeler and Hendon 2004), and is based on the multivariate EOF analysis of daily  
62 anomalies of the zonal wind at 850h Pa and outgoing long-wave radiation (OLR). Other MISO  
63 indices (Kikuchi et al. 2012; Goswami et al. 1999) are based on similar EOF and EEOF tech-  
64 niques. These covariance-based approaches in general capture the spatiotemporal MISO patterns  
65 reasonably well and isolate the northeastward-propagating intraseasonal periodicity band from

66 high-frequency band (Suhas et al. 2013; Abhilash et al. 2014a,b). Yet, the seasonal extraction and  
67 longitudinal averaging in computing these indices are sometimes ad hoc and can potentially lead  
68 to loss of predictive information or mixing with other modes. In addition, these covariance-based  
69 techniques have potential inadequacy in capturing the rare/extreme events in complex nonlinear  
70 dynamics (Crommelin and Majda 2004) which have significant societal and economic impacts.

71 Recently (Sabeerali et al. 2017), a new MISO index based on the Nonlinear Laplacian Spectral  
72 Analysis (NLSA) (Giannakis and Majda 2012b,a) technique was developed. NLSA is a non-  
73 linear data analysis technique that combines ideas from lagged embedding (Packard et al. 1980;  
74 Sauer et al. 1991), machine learning (Coifman and Lafon 2006; Belkin and Niyogi 2003), adap-  
75 tive weights and spectral entropy criteria to extract spatiotemporal modes of variability from high-  
76 dimensional time series. These modes are computed utilizing the eigenfunctions of a discrete  
77 Laplace-Beltrami operator, which can be thought of as a local analog of the temporal covariance  
78 matrix employed in EOF and EEOF techniques, but adapted to the nonlinear geometry of data gen-  
79 erated by complex dynamical systems. A key advantage of NLSA over classical covariance-based  
80 techniques is that NLSA by design requires no ad hoc detrending or spatiotemporal filtering of the  
81 full data set and captures both intermittency and low frequency variability (Giannakis and Majda  
82 2012a,b, 2013; Giannakis et al. 2012). Therefore, the NLSA-based MISO index provides an ob-  
83 jective identification of the MISO patterns in noisy precipitation data. In addition, as reported in  
84 Sabeerali et al. (2017), the NLSA MISO modes have higher memory and predictability, stronger  
85 amplitude and higher fractional explained variance over the western Pacific, Western Ghats, and  
86 adjoining Arabian Sea regions, and more realistic representation of the regional heat sources over  
87 the Indian and Pacific Oceans compared with those extracted via EEOF analysis. Other applica-  
88 tions of NLSA beyond the capability of EOF and EEOF in capturing both the intermittent and

89 low-frequent modes in climate, atmosphere and ocean can be found in Székely et al. (2016a,b);  
90 Slawinska and Giannakis (2016); Giannakis and Majda (2012a, 2011); Brenowitz et al. (2016).

91 In this article, we assess the predictability of the MISO as measured through precipitation. A  
92 recent systematic strategy for data driven physics constrained low-order stochastic modeling of  
93 time series (Majda and Harlim 2013; Harlim et al. 2014) is applied to the two-dimensional MISO  
94 indices from NLSA (Sabeerali et al. 2017). The result is a four dimensional nonlinear stochastic  
95 model for the two MISO variables and two hidden variables. This low-order model involves  
96 correlated multiplicative noise defined through energy conserving nonlinear interactions between  
97 the observed and hidden variables as well as additive stochastic noise. The special structure of  
98 this nonlinear stochastic model allows effective data assimilation algorithm for determining the  
99 initial ensemble of the hidden variables that facilitates the ensemble prediction scheme. This  
100 nonlinear low-order stochastic model has been shown to have significant skill for determining the  
101 predictability limits of the large-scale cloud patterns of both the boreal winter MJO and boreal  
102 summer intraseasonal oscillations (Chen et al. 2014; Chen and Majda 2015a) as well as improving  
103 the prediction skill of the RMM indices (Chen and Majda 2015b). Then with the predicted MISO  
104 indices in hand, an effective and practical spatiotemporal reconstruction algorithm is developed,  
105 which overcomes the fundamental difficulty in most data decomposition techniques with lagged  
106 embedding that require extra information beyond the predicted time series in the future.

107 The remainder of this article is organized as follows. Section 2 describes the precipitation dataset  
108 and the MISO indices obtained from the NLSA technique. Section 3 presents the physics con-  
109 strained low-order nonlinear stochastic model as well as the calibration and the effective prediction  
110 algorithm. The results of predicting the MISO indices are reported in Section 4 and the predic-  
111 tion of the spatiotemporal reconstructed patterns are shown in Section 5. Section 6 discusses the  
112 possibility of shortening the training period to only 3 years and then illustrates the discrepancy

113 in forming the MISO indices with different lagged embedding sizes. Summary conclusions are  
114 included in Section 7.

## 115 **2. The Precipitation MISO Indices from NLSA**

116 The dataset utilized here is the daily Global Precipitation Climatology Project (GPCP) rainfall  
117 data (Huffman et al. 2001) over the Asian summer monsoon region ( $20^{\circ}\text{S}$ - $30^{\circ}\text{N}$ ,  $30^{\circ}\text{E}$ - $140^{\circ}\text{E}$ ) for  
118 period 1997-2014. The spatial resolution of this dataset is  $1^{\circ}\times 1^{\circ}$ , amounting to  $d = 5500$  grid  
119 points for the Asian summer monsoon region.

120 NLSA is applied to the daily GPCP dataset with a lagged embedding window of  $q = 64$  days, an  
121 ideal choice for the intraseasonal time scale. A variety of extended spatial precipitation patterns  
122 emerge from the analysis but the focus here is on the two spatial patterns associated with MISO  
123 with time series depicted in Figure 1. The details of applying NLSA to daily GPCP dataset have  
124 already been described in Sabeerali et al. (2017) and are thus omitted here. It is evident from  
125 Figure 1 that these patterns are active in boreal summer and quiescent in boreal winter. It was  
126 shown in Sabeerali et al. (2017) that the NLSA MISO modes display the characteristic pattern of  
127 northeastward propagating anomalies associated with the MISO. A case study there also revealed  
128 three consecutive MISO events in the NLSA MISO modes in the boreal summer of 2004, the  
129 onset and demise phases of which are highly consistent with observations. These facts indicate  
130 that the time series depicted in Figure 1 give a reasonable representation of the full life cycle of the  
131 northward propagating boreal summer convection band and can be utilized to determine the phase  
132 and amplitude of the poleward-propagating rainfall anomalies associated with the MISO. Below,  
133 we utilize the terminology, MISO indices, for the two time series in Figure 1.

### 134 3. The Low-Order Nonlinear Stochastic Model

135 Denote by  $u_1$  and  $u_2$  the two components, MISO 1 and MSIO 2, depicted in Figure 1. The  
 136 probability distribution functions (PDFs) for  $u_1$  and  $u_2$  are highly non-Gaussian with fat tails  
 137 that indicate the temporal intermittency in the large scale precipitation patterns associated with  
 138 the MSIO. The following family of low-order stochastic models are proposed to describe the  
 139 intermittent variability of the time series  $u_1$  and  $u_2$ :

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

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

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

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

140 where

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

141 In (1),  $u_1$  and  $u_2$  are the two observed MISO variables while  $v$  and  $\omega_u$  are hidden unobserved  
 142 variables which represent the stochastic damping and stochastic phase, respectively. In (1),  
 143  $\dot{W}_{u_1}, \dot{W}_{u_2}, \dot{W}_v$  and  $\dot{W}_\omega$  are independent white noise. The time periodic damping in the equations  
 144 in (1a) and (1b) is utilized to crudely model the active summer season and the quiescent winter  
 145 season in the seasonal cycle. The hidden variables  $v, \omega_u$  interact with the observed MISO vari-  
 146 ables  $u_1, u_2$  through energy conserving nonlinear interactions following the systematic physics  
 147 constrained nonlinear regression strategies for time series developed recently (Majda and Harlim  
 148 2013; Harlim et al. 2014). The energy-conserving nonlinearity is easily seen by multiplying (1a)–  
 149 (1d) by  $u_1, u_2, v$  and  $\omega_u$ , respectively, and then these equations sum up. The energy change in the  
 150 quadratic nonlinear terms cancels with each other and thus the energy due to the nonlinear inter-



151 action is conserved. The low-order stochastic nonlinear models in (1) are fundamentally different  
152 from those utilized earlier (Kondrashov et al. 2013; Kravtsov et al. 2005) which allow for nonlinear  
153 interactions only between the observed variables  $u_1, u_2$  and only special linear interactions with  
154 layers of hidden variables. The physics constrained nonlinear low-order stochastic model (1)–(2)  
155 has been shown to have significant skill for determining the predictability limits of the large-scale  
156 cloud patterns of both the boreal winter MJO and boreal summer intraseasonal oscillations (Chen  
157 et al. 2014; Chen and Majda 2015a) as well as improving the prediction skill of the RMM in-  
158 dices by incorporating a new information-theoretic strategy in the training phase (Chen and Majda  
159 2015b).

#### 160 *a. Calibration of the Nonlinear Stochastic Model*

161 The parameters of the stochastic model in (1)–(2) are calibrated by fitting the highly non-  
162 Gaussian PDFs and autocorrelations of the two MISO variables  $u_1, u_2$  in the training period from  
163 1998 to 2007 as shown in Figure 1. Table 1 records the optimal parameter values while Figure 2  
164 displays the skill of the stochastic model with these parameters in recovering the statistics of the  
165 two MISO indices. Panels (a) and (b) show that the stochastic model succeeds in capturing the  
166 autocorrelations almost perfectly for a three-month duration and even the wiggles that appear with  
167 lags around one year. Panel (c) shows that the stochastic model captures the highly non-Gaussian  
168 fat-tailed PDF of the two MISO indices due to intermittency. Panel (d) shows that the power  
169 spectrum of the two MISO indices from the data and those from the stochastic model match very  
170 well. The optimal parameters in the stochastic model from Table 1 have been determined by sys-  
171 tematically minimizing the information distance of the equilibrium PDF of the stochastic model  
172 compared with that of the actual data (Majda and Gershgorin 2010, 2011). Details are present-

173 ed in the Appendix A. Importantly, the model statistics are robust with respect to the parameter  
174 variations around these optimal values (See Appendix A).

175 *b. Prediction Algorithm and Data Assimilation of the Hidden Variables*

176 The ensemble prediction algorithm is applied to the nonlinear low-order stochastic model (1) for  
177 predicting the MISO time series. The algorithm involves running the forecast model (1) forward  
178 in time given the initial values. The initial data of the two state variables  $\mathbf{U} = (u_1, u_2)$  are obtained  
179 directly from the observations, i.e., MISO 1 and MISO 2 indices. The more important and chal-  
180 lenging issue is to determine the initial ensemble of the two hidden variables  $\Gamma = (v, \omega_u)$ . To this  
181 end, an active data assimilation algorithm is incorporated into the ensemble forecasting scheme.

182 The estimates of the hidden parameters  $\Gamma = (v, \omega_u)$  during the training period and initialization  
183 of these parameters during the prediction phase exploit the special structure of the nonlinear low-  
184 order stochastic model (1). The equations in (1) are a conditional Gaussian system with respect  
185 to the observations  $\mathbf{U} = (u_1, u_2)$ , meaning that once  $u_1$  and  $u_2$  are given the time evolution of the  
186 distributions of  $\Gamma = (v, \omega_u)$  is Gaussian. Such special feature of (1) allows the closed analytic  
187 equations for the conditional Gaussian distributions of the hidden parameters  $\Gamma = (v, \omega_u)$  obtained  
188 from the posterior estimations in the Bayesian framework (Liptser and Shiryaev 2001). Appendix  
189 B contains the details and explicit equations. We utilize this fact to construct an initial ensemble  
190 for forecasting at each time instant in both the training and prediction phases for  $t \in [t_0, t_1, \dots, t_s]$   
191 in the following way.

- 192 1. Starting from a “burn in” time  $t_-$  earlier than  $t_0$  with arbitrary initial conditions for  $\Gamma$ , solve the  
193 associated analytic formula (B2) until time  $t_0$  to obtain the conditional Gaussian distribution  
194  $p_0(\Gamma|\mathbf{U}(t_0))$ . The initial ensemble of the hidden variables  $\Gamma = (v, \omega_u)$  for prediction starting  
195 from  $t_0$  is drawn from this distribution.

- 196 2. The initial ensemble for prediction starting from the next time  $t_1$  is drawn from  $p_1(\Gamma|\mathbf{U}(t_1))$ ,  
 197 where  $p_1(\Gamma|\mathbf{U}(t_1))$  is solved by running the analytic formula (B2) forward from time  $t_0$  to  $t_1$   
 198 with initial value  $p_0(\Gamma|\mathbf{U}(t_0))$ .
- 199 3. Following the same procedure, the initial distributions of the hidden variables  $\Gamma = (v, \omega_u)$  for  
 200 prediction starting from each time  $t_i$  are obtained “on the fly” by running the analytic formula  
 201 (B2) forward from time  $t_{i-1}$  to  $t_i$  with initial value  $p_{i-1}(\Gamma|\mathbf{U}(t_{i-1}))$  when the new observations  
 202 up to  $\mathbf{U}(t_i)$  are available.

203 In the prediction below with (1), we use  $N$  ensemble members with  $N = 50$ .

#### 204 4. Results of Predicting the MISO indices

205 With the optimal parameters from Table 1 and the ensemble initialization scheme described in  
 206 Section 3b, the prediction skill of the stochastic model in (1) for the six year prediction period from  
 207 year 2008 to 2013 is presented here. The skill scores of ensemble mean prediction as a function  
 208 of lead time (days) in different years are shown in Figure 3 and the comparison of the ensemble  
 209 mean prediction and the truth at lead times of 15 and 25 days for MISO 1 index for all six years  
 210 are shown in Figure 4. Here, the skill scores adopted are the root-mean-squared (RMS) error and  
 211 pattern correlation (Corr):

$$\text{RMS error}(\mathbf{U}_t, \mathbf{U}_t^{\text{pred}}) = \sqrt{\frac{\sum_{t=1}^n ((u_{1,t} - u_{1,t}^{\text{pred}})^2 + (u_{2,t} - u_{2,t}^{\text{pred}})^2)}{n}},$$

$$\text{Corr}(\mathbf{U}_t, \mathbf{U}_t^{\text{pred}}) = \frac{\sum_{t=1}^n (u_{1,t} u_{1,t}^{\text{pred}} + u_{2,t} u_{2,t}^{\text{pred}})}{\sqrt{\sum_{t=1}^n (u_{1,t}^2 + u_{2,t}^2)} \sqrt{\sum_{t=1}^n ((u_{1,t}^{\text{pred}})^2 + (u_{2,t}^{\text{pred}})^2)}},$$

212 where  $n$  is the number of the points in the time series, and  $\mathbf{U}_t = (u_{1,t}, u_{2,t})$  and  $\mathbf{U}_t^{\text{pred}} = (u_{1,t}^{\text{pred}}, u_{2,t}^{\text{pred}})$   
 213 are the truth and predicted time series, respectively. It is shown in Figures 3 and 4 that the 15 day  
 214 predictions are very skillful and even the 25 day predictions have highly significant skill in most

215 years. Among different years, year 2010 has useful predictions for about 20 days while year 2011  
216 and 2013 have skillful predictions around 25-30 days. In year 2008, 2009 and 2012, there is a  
217 significant prediction skill out to more than 50 days. Here, useful predictions are defined by 1)  
218 the RMS error in the prediction is less than the standard deviation of the truth at the equilibrium  
219 and 2) the pattern correlation between the predicted signal and the truth is above 0.5. Importantly,  
220 the prediction here yields a significantly higher skill than the conventional EEOF based indices  
221 (Suhas et al. 2013).

222 Both the phase and amplitude of MISO activity play important roles in determining the pre-  
223 diction skill in different years. For example, year 2008 has an overall strong and regular MISO  
224 activity during the whole monsoon season that results in a long predictability, while the signal to  
225 noise ratio in year 2010 is smaller than other years and thus the predictability is greatly affected.  
226 Note that although year 2009 is a drought year with weak MISO activity during the late monsoon  
227 season (September), the MISO activity in other months of 2009 boreal summer remains strong  
228 and the overall prediction skill is high. From the limited sample size (12 years) of our analysis, it is  
229 hard to derive relationship between predictability of MISO and interannual variability of monsoon.  
230 However, it appears that the drought years do not necessarily have low predictability.

231 In addition to the ensemble mean prediction, the ensemble spread that indicates the predictive  
232 uncertainty is another important indicator of the prediction skill. Figure 5 shows the ensemble  
233 predictions including the ensemble spread for the six years, beginning at three different dates:  
234 April 1, June 1 and October 1. It is clear from Figure 1 that April 1 is a time at the transition  
235 between the quiescent phase and the active phase of the MISO indices; June 1 is a starting date in  
236 the active mature phase while October 1 is a starting date in the decaying phase of MISO activity.  
237 As shown in Figure 5, the ensemble mean predictions for the April 1 starting date do not have any  
238 long range skill but the ensemble spread automatically predicts this lack of skill and the envelope

239 of the ensemble predictions contains the true signal for all years and forecast times including the  
240 return to skill in the winter quiescent phase. The forecasts from June 1 obviously have skill from  
241 both the mean and ensemble spread for all years for moderate to long lead times. The forecasts  
242 starting from October 1 have both an accurate mean and small ensemble spread for all six years  
243 and for very long times.

244 It is easy to perform twin prediction experiments with the perfect nonlinear stochastic model  
245 in (1)–(2) where 10 year training segments of the data generated from the model are utilized to  
246 make 6 year forecasts. It is significant that this internal prediction skill of the stochastic model is  
247 comparable to its skill in predicting the MISO indices from observations (not shown here). This  
248 lends support to the fact that the nonlinear stochastic model in (1)–(2) can accurately determine  
249 the predictability limits of the two MISO indices in Figure 1.

## 250 **5. The Spatiotemporal Reconstruction**

251 With the predicted MISO indices in hand, the final step is to recover the spatiotemporal MISO  
252 patterns in physical space. This requires the combination of time series and spatial bases.

### 253 *a. Method*

254 Let  $z_i$  be an  $d$ -dimensional vector of gridded precipitation values over the South Asia monsoon  
255 region at time  $i$ . Here,  $i$  is an integer ranging from 1 to  $n$ , representing the period of training  
256 phase. The first step in NLSA is to construct a higher-dimensional, time lagged embedding dataset  
257 utilizing Takens' method of delay. Denote  $q$  be the lagged embedding window size. Then the

258 lagged embedding matrix can be written as

$$X = \begin{pmatrix} z_1 & z_2 & \cdots & z_{n-2q+1} \\ z_2 & z_3 & \cdots & z_{n-2q+2} \\ \vdots & \vdots & \ddots & \vdots \\ z_{q-1} & z_q & \cdots & z_{n-q-1} \\ z_q & z_{q+1} & \cdots & z_{n-q} \end{pmatrix} = \begin{pmatrix} z_1 & z_2 & \cdots & z_{N-q+1} & z_{N-q+2} & \cdots & z_{N-1} & \boxed{z_N} \\ z_2 & z_3 & \cdots & z_{N-q+2} & z_{N-q+3} & \cdots & \boxed{z_N} & z_{N+1} \\ \vdots & \vdots & \ddots & \vdots & \cdots & \cdots & \cdots & \vdots \\ z_{q-1} & z_q & \cdots & z_{N-1} & \boxed{z_N} & \cdots & z_{N+q-3} & z_{N+q-2} \\ z_q & z_{q+1} & \cdots & \boxed{z_N} & z_{N+1} & \cdots & z_{N+q-2} & z_{N+q-1} \end{pmatrix}, \quad (3)$$

259 where  $N = n - 2q + 1$ . Note that  $q$  in (3) is actually  $dq$  but  $d$  is omitted here for notation simplicity.

260 Although the lagged embedding matrix  $X$  in (3) is formed by the raw observational data, the matrix  
 261 associated with each eigenmode, such as the annual mode, semi-annual mode and MISO mode,  
 262 after the spatiotemporal reconstruction has essentially the same structure, except that the  $q$  entries  
 263 that are represented by the same  $z_i$  in (3) may have different values. Therefore, averaging over  
 264 these  $q$  entries finalizes the reconstruction of  $z_i$ . From now on,  $X$  represents the lagged embedding  
 265 matrix containing only the MISO mode. Note that, since  $z_1, \dots, z_{q-1}$  and  $z_{N+1}, \dots, z_{N+q-1}$  appear  
 266 less than  $q$  times in (3), recovering these components requires a longer observational period.

267 The relationship between the spatiotemporal representation  $X$ , the time series (i.e., indices)  $\Phi$   
 268 and the spatial basis  $A$  is simply given by

$$X = A\Phi^T, \quad (4)$$

269 where  $\cdot^T$  stands for the transpose. Clearly, with the predicted MISO indices in hand, the recon-  
 270 struction (4) is easily achieved as long as  $A$  is able to be reached in the predicted phase. Different  
 271 from the 2-Dimensional time series that is predicted by the low-order models, it is a challenge  
 272 to describe and predict the exact evolution of the high-dimensional spatial basis  $A$ . Therefore,  
 273 approximations are typically included in developing  $A$  in the predicted phases. Below, the spatial  
 274 basis  $A$  utilized for prediction is assumed to be a constant matrix that is completely determined in  
 275 the training phase. Note that the stationary assumption of  $A$  is in general not necessary and may

276 even result in some errors in the spatiotemporal reconstruction for nonlinear data. Nevertheless,  
 277 adopting a constant matrix  $A$  greatly reduces the computational cost and facilitates an effective  
 278 and practical prediction algorithm as will be demonstrated soon. As will be shown at the end of  
 279 this section, the approximated reconstruction with such a constant spatial basis is actually highly  
 280 consistent with the truth. In Section 6d, one alternative of creating a non-stationary spatial basis  
 281 will be briefly discussed and compared with the method proposed in this section.

282 In light of both  $X$  and  $\Phi$  in the training period, the spatial pattern  $A$  according to (4) is given by

$$A = cX\Phi, \quad \text{with } c = \frac{1}{\|\Phi\|^2}. \quad (5)$$

283 Then a natural way of performing the spatiotemporal reconstruction of the predicted MISO pat-  
 284 terns is to multiply  $A$  obtained from (5) by the predicted MISO indices. Denote  $X^f$  and  $\Phi^f$  the  
 285 spatiotemporal pattern and indices in the prediction period. The following relation is reached:

$$A \cdot [\Phi; \Phi^f] = [X; X^f], \quad (6)$$

286 where we ignore the transpose in  $\Phi$  and  $\Phi^f$  for notation simplicity.

287 However, as is shown in Appendix C, the fundamental barrier for the direct method (6) to be-  
 288 come practical is that in order to obtain the spatiotemporal patterns at  $s$  lead days, the prediction  
 289 of the time series up to  $s + q$  days is required (Comeau et al. 2016). For example, in the prediction  
 290 of MISO, reconstructing the spatiotemporal pattern for the next day requires the prediction of time  
 291 series up to the next 65 days! In fact, to reach the last  $z_N$  as shown in (3), the information up to  
 292  $z_{N+q-1}$  is required.

293 One remedy to overcome such difficulty is to switch the extra future information as required in  
 294 the time series  $\Phi$  to that in the spatial basis  $A$  by calculating a “predicted spatial basis”  $\tilde{A}$  in the

295 training phase. This  $\tilde{A}$  is obtained by taking advantage of a new lagged embedding matrix  $\tilde{X}$ ,

$$\tilde{X} = \begin{pmatrix} z_q & z_{q+1} & \cdots & z_{n-q} \\ z_{q+1} & z_{q+2} & \cdots & z_{n-q+1} \\ \vdots & \vdots & \ddots & \vdots \\ z_{2q-1} & z_{2q} & \cdots & z_{n-1} \end{pmatrix}, \quad (7)$$

296 which is just  $q - 1$  units shift forward in time with respect to  $X$  in (3). Similar as  $A$ , a new spatial  
297 pattern  $\tilde{A}$  is formed by

$$\tilde{A} = c\tilde{X}\Phi, \quad \text{with } c = \frac{1}{\|\Phi\|^2}. \quad (8)$$

298 Replacing  $A$  by  $\tilde{A}$  in (6) leads to

$$\tilde{X} \cdot [\Phi; \Phi^f] = [\tilde{X}, \tilde{X}^f]. \quad (9)$$

299 As shown in Appendix C, with (9), reconstructing the spatiotemporal patterns at  $s$  days in the future  
300 requires only the prediction of the time series for  $s$  days. Since  $\tilde{A}$  is completely determined in the  
301 training period that involves only straightforward calculations, the formula in (9) is an effective  
302 and practical method for predicting the spatiotemporal patterns.

### 303 *b. Prediction of the Spatial-Temporal Reconstructed Precipitation Fields*

304 Figure 6 includes three phase diagrams of the MISO indices, each containing a length of one-  
305 month period. The corresponding predictions, starting from the first day of each period and lasting  
306 for 30 days, is also included. Among the three periods, a significant skillful prediction is found for  
307 July 2009 while the prediction skill of June 2008 is moderate. The true signal of June 2013 has  
308 a weak amplitude and the corresponding prediction is far from the truth. Below, the prediction of  
309 spatiotemporal patterns based on the improved method (9) is demonstrated, where the ensemble  
310 mean of prediction is utilized for spatiotemporal reconstruction.



311 The skill scores of the predicted spatiotemporal patterns for each of the three periods are shown  
312 in Figure 7. Consistent with the MISO indices, July 2009 has the highest prediction skill and  
313 the useful prediction lasts for 40 days. On the other hand, a higher pattern correlation is found  
314 in predicting the spatiotemporal patterns of June 2008 than that of June 2013, where the useful  
315 prediction of June 2008 is up to around 22 days. Note that, different from predicting the MISO  
316 indices, the skill scores of predicting the spatiotemporal pattern do not decrease monotonically  
317 as a function of lead time. This is due to the stationary approximation of the spatial basis  $\tilde{A}$   
318 in the prediction period. In addition, this averaged spatial basis may also lead to the amplitude  
319 underestimation in the predicted spatiotemporal reconstruction. A direct remedy is to compute the  
320 ratio of  $\|A\|$  and  $\|\tilde{A}\|$  in the training period and multiply this constant ratio in prediction. In fact,  
321 the value of  $\|\tilde{A}\|$  decreases with the increases in  $q$  value. This is because the correlation between  
322 the spatial basis and the time series that with a phase lag becomes weaker when the lag increases.

323 Figures 8 and 9 compare the truth and the predicted spatiotemporal patterns of July 2009 and  
324 June 2008, respectively. The predicted patterns for the whole July 2009 are highly consistent  
325 with the truth, especially in the regions of Indian subcontinent and Bay of Bengal. On the other  
326 hand, despite the skillful prediction up to 20 days lead time, significant errors in the spatiotem-  
327 poral patterns appear for longer time predictions of June 2008 due to the failure in predicting the  
328 precipitation in regions such as the Indian Ocean.

329 To understand the approximated error in  $\tilde{A}$ , the true spatiotemporal patterns from NLSA, which  
330 is validated in Sabeerali et al. (2017), and the approximated patterns based on  $\tilde{A}\Phi$  for July 2009  
331 are compared in the first two rows of Figure 10. Note that the truth of  $\Phi$  is adopted here to exclude  
332 the error in the prediction of the time series. Despite the stationary approximation in creating the  
333 spatial basis  $\tilde{A}$ , the time series  $\Phi$  from NLSA remains highly nonlinear and intermittent. Clearly,

334 the approximated patterns are remarkably consistent with the truth, where the non-Gaussian and  
335 intermittent features in the spatiotemporal patterns are both retained to a high extent.

## 336 **6. Discussions**

### 337 *a. Prediction with a 3-year short training period*

338 A typical situation in climate science is that only a short period of observational data is avail-  
339 able. This actually leads to one of the fundamental difficulties in prediction utilizing most non-  
340 parametric methods that require a huge amount of data for training. Suitable models that are  
341 able to describe essential characteristics of the data are usually preferred since they allow a much  
342 shorter training period. Recall in previous sections, 10 years of observations (1998-2007) were  
343 adopted for model calibration and the prediction skill were assessed for the remaining 6 years  
344 (2008-2013). Although this 10-year training window is already much shorter than that required  
345 by most non-parametric methods, it is important to understand whether an even shorter training  
346 period is possible here for the nonlinear model to obtain the information in nature.

347 To this end, a very short training period involving only the first three years of the time series  
348 (1998-2000) is adopted here for model calibration. Figure 11 compares the statistics of the MISO  
349 time series with different lengths, including this short 3-year training period (1998-2000), the  
350 10-year training period adopted in previous sections (1998-2007) and the full MISO period (1998-  
351 2013). The fact that the statistics of the 10-year training period and the full MISO period almost  
352 perfectly match each other indicates the sufficiency of the 10-year training period in obtaining  
353 the unbiased information. On the other hand, the 3-year training period, including one weak  
354 year (1998), one moderate year (1999) and one strong year (2000) of MISO activity, also has  
355 highly consistent statistics with those associated with the full MISO time series, including the

356 non-Gaussian fat-tailed PDFs, the power spectrums and the autocorrelations up to 1.5 months.  
357 Therefore, the information of the full MISO indices are well reflected in this short 3-year training  
358 period. Due to the robustness of the model parameters (Appendix A), the calibrated parameters  
359 based on this 3-year short training period are nearly the same as the optimal parameters shown  
360 in Table 1. Importantly, this short training period allows the study of prediction skill for a long  
361 period back to year 2001 and the results are roughly reported here.

362 Figure 12 shows the skill scores and the predicted signals based on the ensemble mean prediction  
363 from year 2001 to 2007, analogous to those in Figure 3 and 4 from year 2008 to 2013. The useful  
364 prediction of these 7 years all exceeds 25 days, where in particular the skillful predictions in year  
365 2001, 2003 and 2007 are more than 40 days. Among these 7 years, year 2002 and 2004 are  
366 recorded as drought years. A significant error is found in predicting the subdued MISO activity  
367 during August and September of year 2002, which explains its lower overall prediction skill than  
368 most of the other years. On the other hand, despite being a drought year, MISO activity during  
369 2004 is persistently strong throughout the boreal summer. The major error in predicting MISO  
370 indices of year 2004 is in fact due to the model's failure in capturing the extremely slow oscillation  
371 frequency during August and September.

372 We have also check the model statistics and prediction skill by utilizing any three consecutive  
373 years between 1998-2013 as the training phase. Despite the discrepancy in the signal variance  
374 due to the strength of MISO activity in different years, the fat tails in the non-Gaussian PDFs,  
375 the peak of the power spectrums and the autorrelations up to 1.5 months all resemble those of  
376 the full MISO time series. Importantly, the ensemble prediction skill does not have significant  
377 deterioration based on different training periods.

378 *b. MISO indices based on different lagged embedding window sizes and the corresponding pre-*  
379 *diction skill*

380 Recall that the two MISO indices shown in Figure 1 and studied throughout this article were  
381 obtained by applying NLSA to the precipitation data with a lagged embedding window of length  
382  $q = 64$  days. Adopting  $q = 64$  is natural since it is an ideal choice for representing the intraseasonal  
383 time scale and such lagged embedding window size was utilized for defining the large-scale cloud  
384 patterns of the MJO and monsoon in previous works (Chen et al. 2014; Chen and Majda 2015b,a;  
385 Tung et al. 2014). On the other hand, EEOF was also widely utilized in defining the MISO indices  
386 in literature (Suhas et al. 2013; Kikuchi et al. 2012), which involves removing the climatological  
387 mean, first a few harmonics of the seasonal cycle and then applying a much shorter embedding  
388 window with 15-20 days. Therefore, it is important to study the difference in the MISO indices by  
389 applying NLSA with different lagged embedding window sizes.

390 Figure 13 shows the resulting MISO indices by applying NLSA with  $q = 64, 48$  and  $34$  as well  
391 as the corresponding statistics. Different from  $q = 64$ , the MISO indices with  $q = 48$  and  $34$   
392 have active phases in both boreal summer and winter, implying that the obtained MISO indices  
393 contain the components of the boreal winter MJO, and the associated PDFs are nearly Gaussian.  
394 In addition to being polluted by the boreal winter signal, these time series, especially with  $q = 34$ ,  
395 also contain bi-annual, annual and semi-annual cycles, as indicated by the large bursts in the low-  
396 frequent band of the power spectrum. Another significant discrepancy with different  $q$  values  
397 lies in the causality between the two components of the MISO indices. With  $q = 64$ , the cross-  
398 correlation functions have significant peaks at lags around 12 days, which is nearly 1/4 of the  
399 averaged oscillation frequency and indicates the quadrature structure of MISO 1 and MISO 2. On  
400 the other hand, the cross-correlations  $R_{12}(t)$  and  $R_{21}(t)$  with  $q = 48$  and  $q = 34$  remain close to

401 zero, and the maximum value of the lagged correlation between MISO 1 and MISO 2 indices is  
402 less than 0.3 (not shown here). These facts imply that MISO 1 and MISO 2 are nearly uncorrelated  
403 and therefore model errors appear in fitting the cross-correlations utilizing the nonlinear low-order  
404 model (1). Finally, the fast decay of autocorrelations  $R_{11}(t)$  and  $R_{22}(t)$  with  $q = 48$  and 34 implies  
405 deterioration in the predictability of the MISO indices.

406 Figure 14 shows the prediction skill with different  $q$ . Here useful prediction is defined in the  
407 same way as that in Section 4: 1) the RMS error in the prediction is less than the standard deviation  
408 of the truth at the equilibrium and 2) the pattern correlation between the predicted signal and the  
409 truth is above 0.5. In addition to illustrating the prediction skill for the whole year, the prediction  
410 skill conditioned on the boreal summer time (June to September) is also emphasized. As expected,  
411 with the decrease in  $q$ , the overall prediction skill becomes worse. Nevertheless, conditioned on  
412 the boreal summer time, the prediction with  $q = 48$  remains quite skillful and in particular the  
413 15-day lead time prediction is highly consistent with the truth. This is, however, not true for the  
414 prediction with  $q = 34$ , where the useful prediction only lasts for 10-12 days in terms of both the  
415 whole year and only the boreal summer time.

416 *c. Significant prediction skill of the precipitation MISO indices with parameters calibrated from*  
417 *OLR dataset*

418 Most tropical rainfall is convective, which implies that OLR, a proxy for the convection, is a  
419 potential candidate to describe the precipitation in tropics. Positive (negative) OLR anomalies are  
420 associated with reduced (increased) cloudiness, hence suppressed (enhanced) deep convection.  
421 Due to the possible relationship between the OLR and the tropical precipitation anomalies, it is  
422 important to understand the skill of the low-order nonlinear stochastic model (1)–(2) in predicting  
423 the MISO indices with parameters calibrated from OLR dataset.

424 In Chen and Majda (2015b), the low-order nonlinear stochastic model (1)–(2) was adopted to  
425 predict the two boreal summer intraseasonal oscillation (BSISO) modes obtained by applying  
426 NLSA to the brightness temperature, a highly correlated variable with OLR, within the equato-  
427 rial tropical belt from  $15^{\circ}S$  to  $30^{\circ}N$ . The dataset utilized there was Cloud Archive User Service  
428 (CLAUS) Version 4.7. To explore the strength of the correlation between OLR and precipitation,  
429 the parameters in Chen and Majda (2015b) are applied to (1)–(2) to calibrate and predict the pre-  
430 cipitation MISO indices. For simplicity, these parameters are named as OLR-based parameters.

431 The OLR-based parameters are listed in the second row of Table 1 with two minor modifications.  
432 First, since the time series in Chen and Majda (2015b) were started from September instead of  
433 January, the phase parameter  $\phi$  in Chen and Majda (2015b) is modified accordingly. Second,  
434 due to the general negative correlation between OLR and precipitation, the sign of the oscillation  
435 frequency  $a$  in Chen and Majda (2015b) is flipped. In fact, as shown in Table 1, the OLR-based  
436 parameters are quite similar to the optimal parameters utilized in the previous sections.

437 Panels (a)–(d) of Figure 15 show the model statistics with the OLR-based parameters. The  
438 autocorrelations, power spectrums and non-Gaussian fat-tailed PDFs are all quite consistent with  
439 the truth, implies nearly identical statistical and dynamical features in describing the precipitation  
440 MISO with the OLR-based parameters. The slight underestimation of the variance with the OLR-  
441 based parameters is due to the gap in the units of the two variables, which can easily be remedied  
442 by applying an amplitude rescaling to the OLR variable and is a secondary issue here. Panels  
443 (e) and (f) compare the 25-day lead ensemble mean predictions using the optimal parameters and  
444 the OLR-based parameters of three different years from 2008 to 2010 with strong, moderate and  
445 weak MISO activities, respectively. In fact, except a light underestimation of the amplitude due  
446 to the insufficient amplitude of the seasonal cycle damping  $f_i$ , the prediction utilizing OLR-based  
447 parameters has almost the same skill as that utilizing the optimal parameters in all the three years.

448 The results from both model calibration and prediction confirm a strong (negative) correlation  
449 between OLR and precipitation anomalies (Sabeerali et al. 2017).

450 *d. Defects of creating the spatial basis based on the running average over the raw data*

451 Recall in Section 5, despite a stationary spatial basis  $\tilde{A}$  being utilized, the approximated spa-  
452 tiotemporal reconstructed pattern  $\tilde{A}\Phi$  is highly consistent with nature (Figure 10). In addition to  
453 the strong nonlinear and intermittent time series  $\Phi$ , adopting the spatiotemporal pattern  $\tilde{X}$  from  
454 NLSA in determining  $\tilde{A}$  is another important reason for such high consistency, which will be  
455 emphasized in this section.

456 Note that the true spatiotemporal patterns based on the traditional linear methods, such as EOF  
457 and EEOF, are typically reached by applying bandpass filter or running average on the raw obser-  
458 vations over a certain prescribed window. Below, running average is applied to the raw data, the  
459 results of which in the training period are then utilized to form the spatial basis. In the prediction  
460 period, such spatial basis is multiplied by the NLSA MISO indices for the spatiotemporal recon-  
461 struction. Since the same NLSA MISO indices are utilized here and in Section 5, the discrepancy  
462 in the reconstructed spatiotemporal patterns completely lies in the spatial basis.

463 The details of the spatiotemporal reconstruction mentioned above is presented below. Note that  
464 a non-stationary spatial basis is adopted here based on phase decomposition in the training period.  
465 First, a  $q$  day running average is applied on raw rainfall datasets and the daily rainfall anomaly  
466 is computed at each grid point in the training period. Then the whole phase space of MISO 1  
467 and MISO 2 is divided into  $S$  disjoint pieces, named as phases. The spatial pattern of each phase  
468 is computed by conditional averaging of this rainfall anomalies subject to the criteria that the  
469 instantaneous MISO indices have amplitude being greater than 1 standard deviation and belong to  
470 the corresponding phase. In the prediction stage, multiplying the magnitude of the predicted MISO

471 indices by the selected spatial basis according to the location of the predicted MISO indices in the  
472 phase space. This results in the spatiotemporal reconstructed pattern. A typical situation involves  
473 dividing the whole phase space into  $S = 8$  pieces with equal areas (See Figure 6) as was done  
474 in RMM and other MJO and monsoon indices (Wheeler and Hendon 2004; Székely et al. 2016a;  
475 Suhas et al. 2013; Lee et al. 2013). Since all the phases and the associated average spatial bases are  
476 determined in training period, this spatiotemporal prediction is also practical and computational  
477 efficient.

478 We apply this phase decomposition method with  $S = 8$  to reach the spatial-temporal patterns  
479 of July 2009 and the results are shown in the third row of Figure 10. Significant differences are  
480 found compared with the truth (first row). Particularly, a reversed drought/flood pattern from July  
481 1 to July 11 in the India subcontinent based on these two methods is observed and the amount of  
482 precipitation in Indian Ocean and Arabian Sea is significantly different within the whole month.  
483 Such errors result from the spatial basis that is determined by applying the running average over  
484 the raw data. In fact, as was pointed out in Sabeerali et al. (2017) that many important MISO  
485 features are not well represented by those linear methods, including underestimating the fractional  
486 variance over Western Ghats and the failure in capturing variability over Indo-West Pacific region  
487 which is particularly crucial in determining the propagation characteristics of MISO (Sabeerali  
488 et al. 2017). This indicates the importance of the NLSA spatial patterns in addition to the nonlinear  
489 and intermittent time series.

490 Replacing the non-stationary spatial basis resulting from the modes based on the running average  
491 of the raw data by the NLSA modes is a potential way to improve the spatiotemporal reconstruction  
492 in the prediction stage. Yet, the phase decomposition method has an obvious drawback that the  
493 predicted spatiotemporal patterns are discontinuous in time when the corresponding spatial basis  
494 transits from one phase to another. One remedy to the discontinuity issue is to introduce a smooth



495 transition between different phases such as adopting a convolution with a Gaussian kernel. This  
496 remains as a future work.

## 497 **7. Conclusions**

498 A recently developed technique for nonlinear time series analysis NLSA (Giannakis and Majda  
499 2012a,b, 2013) has been utilized to define two MISO indices from the daily GPCP rainfall data  
500 set without detrending or spatiotemporal filtering (Sabeerali et al. 2017). The observed time series  
501 have non-Gaussian fat-tailed PDFs, which is a consequence of intermittency.

502 Systematic strategies for physics constrained regression models (Majda and Harlim 2013; Har-  
503 lim et al. 2014) suggest a four dimensional stochastic model with two hidden variables repre-  
504 senting stochastic damping and random phasing with energy conserving nonlinear feedback in-  
505 teraction (Section 3). In a calibration phase, these models can successfully capture the observed  
506 non-Gaussian PDFs, power spectrums and autocorrelations. The models have a special structure  
507 that allows efficient data assimilation and ensemble initialization algorithms for the hidden vari-  
508 ables. It is shown in Section 4 that the low-order nonlinear stochastic model has been applied to  
509 prediction of the NLSA MISO indices with forecasting skill ranging from 20 to 50 days in dif-  
510 ferent years. Furthermore, the ensemble spread in the stochastic model has been shown to be an  
511 accurate predictive indicator of forecast uncertainty at long range.

512 An effective and practical spatiotemporal reconstruction algorithm is then proposed in Section 5,  
513 which overcomes the fundamental difficulty in most data decomposition techniques with lagged  
514 embedding that require extra information beyond the predicted time series in the future. The  
515 prediction skill of the reconstruction spatiotemporal patterns is consistent with that of the MISO  
516 indices.

517 A few issues are addressed in Section 6. First, the model calibration and prediction with a  
518 3-year short training period is studied. The resulting statistics and prediction skill do not have sig-  
519 nificant deterioration compared with those based on a 10-year training period. This suggests the  
520 advantage of utilizing the low-order nonlinear model (1) over most non-parametric methods in pre-  
521 dicting the MISO indices from a practical point of view. Second, the NLSA MISO indices based  
522 on different lagged embedding window sizes are compared. The resulting MISO indices with a  
523 shorter lagged embedding window size ( $q = 48$  and  $q = 34$ ) are polluted with other variabilities  
524 and the corresponding overall prediction skill is greatly affected. Nevertheless, with  $q = 48$  days  
525 lag, the prediction conditioned on the boreal summer time remains still skillful. Thirdly, the low-  
526 order nonlinear stochastic model with OLR-based parameters remains high skill in both fitting  
527 the non-Gaussian statistics and predicting the precipitation MISO indices, implying a significant  
528 correlation between the tropical precipitation and OLR. Finally, algorithms with potential abilities  
529 to improve the prediction skill of the spatiotemporal patterns are briefly discussed, the implemen-  
530 tation of which remains as future works.

531 Developing more effective and accurate spatiotemporal reconstruction algorithm remains as one  
532 of the future works. In fact, clustering method is a promising technique for recovering more  
533 detailed features of spatial basis conditioned on different phases.

534 *Acknowledgments.* The research of A.J.M is partially supported by the Office of Naval Research  
535 Grant ONR MURI N00014-16-1-2161 and the New York University Abu Dhabi Research Institute.  
536 N.C. is supported as a postdoctoral fellow following A.J.M's ONR MURI Grant. C.T.S, R.S.A and  
537 A.J.M also acknowledge the support from Monsoon Mission of the Ministry of Earth Sciences  
538 (MoES), Government of India (Grant No. MM/SERP/NYU/2014/SSC-01/002). The research of

539 C. T. S and R.S. A is also supported by the New York University Abu Dhabi Research Institute.  
540 The authors thank Dimitrios Giannakis for useful discussions.

## 541 APPENDIX A

### 542 Calibration of the nonlinear stochastic model with information theory

543 The optimal parameters in the nonlinear low-order stochastic model (1) are calibrated by sys-  
544 tematically minimizing the information distance, i.e., the model error, in the PDF of the model  
545  $\pi^M$  compared with that of the MISO index  $\pi$  (Majda and Gershgorin 2010, 2011; Kleeman 2002;  
546 Majda and Branicki 2012; Branicki et al. 2013),

$$\mathcal{D}(\pi, \pi^M) = \int \pi \log \left( \frac{\pi}{\pi^M} \right). \quad (\text{A1})$$

547 The model error dependence on the variation of different parameters is shown in Figure A1, which  
548 indicates that the nonlinear low-order stochastic model (1) is robust with respect to the parameters  
549 around their optimal values. The huge model error with the underestimation of  $\sigma_u, f_t$  and  $\gamma$  and  
550 the overestimation of  $d_u$  is due to the failure of capturing the intermittency. Note that the model  
551 error has only a weak dependence on the background phase  $a$  since the contribution of the oscil-  
552 lation in the signal has been averaged out in the time-averaged PDF. However, the parameter  $a$  is  
553 crucial in describing the frequency of intraseasonal oscillation, and it is calibrated by matching the  
554 autocorrelation functions associated with the model and the truth. The other parameters  $d_v, \sigma_v, d_\omega$   
555 and  $\sigma_\omega$  in describing the stochastic processes affect not only on the model error but more on the  
556 autocorrelations and power spectrums as well. A large discrepancy appears in the statistics if these  
557 parameters are outside the optimal range. The parameter  $f_0$  is not an independent parameter given  
558  $d_u$  and  $\gamma$  and therefore we fix its value. The frequency  $\omega_f$  in the time-periodic damping  $v_f(t)$  is  
559 prescribed to be  $2\pi/12$  such that one time unit of the model corresponds to one month in reality.

560 The phase  $\phi$  in  $v_f(t)$  is tuned to make the strong intermittency occur in the boreal summer in  
561 accordance with the MISO indices. Note that none of the parameters is redundant in the nonlin-  
562 ear stochastic model (1). In fact, without the hidden variables  $v$  and  $\omega_u$ , even if the time-period  
563 damping  $v_f(t)$  is able to crudely describe the active phase of BSISO in the reduced linear model,  
564 a distinguished disparity is observed in the model statistics compared with the truth, indicating the  
565 intrinsic barrier (Majda and Gershgorin 2011; Majda and Branicki 2012).

566 Prediction with random suboptimal parameters is also studied. Here the suboptimal parameters  
567 are taken randomly between the two dotted lines in each panel of Figure 7. Comparable prediction  
568 skill is found with these random suboptimal parameters as the optimal parameters.

## 569 APPENDIX B

### 570 Mathematical details of effective data assimilation and prediction algorithm

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

$$d\mathbf{U}_t = [\mathbf{A}_0(t, \mathbf{U}) + \mathbf{A}_1(t, \mathbf{U})\Gamma_t]dt + \Sigma_U(t, \mathbf{U})d\mathbf{W}_U(t), \quad (\text{B1a})$$

$$d\Gamma_t = [\mathbf{a}_0(t, \mathbf{U}) + \mathbf{a}_1(t, \mathbf{U})\Gamma_t]dt + \Sigma_\Gamma(t, \mathbf{U})d\mathbf{W}_\Gamma(t), \quad (\text{B1b})$$

573 where

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

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

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

574 The model (B1) is a conditional Gaussian system conditioned on the observations  $\mathbf{U}$ , meaning that  
 575 once the observations  $\mathbf{U}$  are given the dynamics of  $\Gamma$  in (B1) becomes a Gaussian system (Chen  
 576 and Majda 2016). The special structure of system (B1) allows the closed analytic formulae for the  
 577 evolution of the conditional Gaussian distributions of the hidden parameters  $v$  and  $\omega_u$  (Liptser and  
 578 Shiryaev 2001) obtained in the Bayesian framework:

$$\begin{aligned}
 d\mu_t &= [\mathbf{a}_0(t, \mathbf{U}) + \mathbf{a}_1(t, \mathbf{U})\mu_t]dt + (R_t \mathbf{A}_1^*(t, \mathbf{U}))(\Sigma_U \Sigma_U^*)^{-1}(t, \mathbf{U}) \times \\
 &\quad [d\mathbf{U}_t - (\mathbf{A}_0(t, \mathbf{U}) + \mathbf{A}_1(t, \mathbf{U})\mu_t)dt], \\
 dR_t &= \{ \mathbf{a}_1(t, \mathbf{U})R_t + R_t \mathbf{a}_1^*(t, \mathbf{U}) + (\Sigma_\Gamma \Sigma_\Gamma^*)(t, \mathbf{U}) \\
 &\quad - (R_t \mathbf{A}_1^*(t, \mathbf{U}))(\Sigma_U \Sigma_U^*)^{-1}(t, \mathbf{U})(R_t \mathbf{A}_1^*(t, \mathbf{U}))^* \} dt,
 \end{aligned} \tag{B2}$$

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

581 As a remark, the formulae (B2) are optimal if and only if the signal is generated from system  
 582 (B1). Since our observed signal, i.e., the MISO indices, are not from the nonlinear low-order s-  
 583 tochastic model (B1), the evolutions of the conditional Gaussian distributions (B2) are suboptimal.

## 584 APPENDIX C

### 585 **Details of the spatiotemporal reconstruction**

586 Recall the relation between the spatial basis  $\Phi$ , the time series  $A$  and the spatiotemporal patterns  
 587  $X$  given by the direct method (6)

$$A \cdot [\Phi; \Phi^f] = [X; X^f]. \tag{C1}$$

588 The left and right hand side of (C1) are given respectively by

$$\begin{aligned}
 A \cdot [\Phi; (\Phi^f)] &= \begin{pmatrix} A_1 \\ A_2 \\ \vdots \\ A_q \end{pmatrix} \cdot \left( \Phi_1, \dots, \Phi_N, \Phi_1^f, \Phi_2^f, \dots, \Phi_q^f \right) \\
 &= \begin{pmatrix} A_1\Phi_1 & \cdots & A_1\Phi_N & A_1\Phi_1^f & A_1\Phi_2^f & \cdots & \boxed{A_1\Phi_q^f} \\ A_2\Phi_1 & \cdots & A_2\Phi_N & A_2\Phi_1^f & A_2\Phi_2^f & \cdots & A_2\Phi_q^f \\ \vdots & \cdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ A_{q-1}\Phi_1 & \cdots & A_{q-1}\Phi_N & A_{q-1}\Phi_1^f & \boxed{A_{q-1}\Phi_2^f} & \cdots & A_{q-1}\Phi_q^f \\ A_q\Phi_1 & \cdots & A_q\Phi_N & \boxed{A_q\Phi_1^f} & A_q\Phi_2^f & \cdots & A_q\Phi_q^f \end{pmatrix} \quad (C2)
 \end{aligned}$$

589 and

$$[X; X^f] = \begin{pmatrix} z_1 & z_2 & \cdots & z_{n-2q+1} & z_{n-2q+2} & z_{n-2q+3} & \cdots & z_{n-q} & \boxed{z_1^f} \\ z_2 & z_3 & \cdots & z_{n-2q+2} & z_{n-2q+3} & z_{n-2q+4} & \cdots & \boxed{z_1^f} & z_2^f \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ z_{q-1} & z_q & \cdots & z_{n-q+1} & z_{n-q} & \boxed{z_1^f} & \cdots & z_{q-1}^f & z_{q-1}^f \\ z_q & z_{q+1} & \cdots & z_{n-q} & \boxed{z_1^f} & z_2^f & \cdots & z_q^f & z_q^f \end{pmatrix}, \quad (C3)$$

590 However, the entries with boxes in (C2) and (C3) implies that in order to obtain the spatiotem-  
 591 poral pattern at one lead time unit, the time series up to  $q$  lead time units are required. In other  
 592 words, predicting  $q = 64$  days of the time series in the future is only sufficient to achieve the  
 593 spatiotemporal pattern for 1 day forward. Therefore, this method is not practical.

594 On the other hand, the improved method based on the new spatial basis  $\tilde{A}$  in (8). Recall the  
 595 relationship in (C4),

$$\tilde{A} \cdot [\Phi; \Phi^f] = [\tilde{X}, \tilde{X}^f]. \quad (C4)$$

596 The left and right hand sides can be written down explicitly,

$$\begin{aligned}
 \tilde{A} \cdot [\Phi; \Phi^f] &= \begin{pmatrix} \tilde{A}_1 \\ \tilde{A}_2 \\ \vdots \\ \tilde{A}_q \end{pmatrix} \cdot \left( \Phi_1, \dots, \Phi_{N-q}, \left| \Phi_{N-q+1}, \dots, \Phi_N, \left| \Phi_1^f, \Phi_2^f, \dots, \Phi_q^f \right. \right. \right) \\
 &= \begin{pmatrix} \tilde{A}_1 \Phi_1 & \cdots & \tilde{A}_1 \Phi_{N-q} & \tilde{A}_1 \Phi_{N-q+1} & \tilde{A}_1 \Phi_{N-q+2} & \cdots & \tilde{A}_1 \Phi_N & \boxed{\tilde{A}_1 \Phi_1^f} & \cdots & \tilde{A}_1 \Phi_q^f \\ \tilde{A}_2 \Phi_1 & \cdots & \tilde{A}_2 \Phi_{N-q} & \tilde{A}_2 \Phi_{N-q+1} & \tilde{A}_2 \Phi_{N-q+2} & \cdots & \boxed{\tilde{A}_2 \Phi_N} & \tilde{A}_2 \Phi_1^f & \cdots & \tilde{A}_2 \Phi_q^f \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \tilde{A}_{q-1} \Phi_1 & \cdots & \tilde{A}_{q-1} \Phi_{N-q} & \tilde{A}_{q-1} \Phi_{N-q+1} & \tilde{A}_{q-1} \Phi_{N-q+2} & \cdots & \tilde{A}_{q-1} \Phi_N & \tilde{A}_{q-1} \Phi_1^f & \cdots & \tilde{A}_{q-1} \Phi_q^f \\ \tilde{A}_q \Phi_1 & \cdots & \tilde{A}_q \Phi_{N-q} & \tilde{A}_q \Phi_{N-q+1} & \boxed{\tilde{A}_q \Phi_{N-q+2}} & \cdots & \tilde{A}_q \Phi_N & \tilde{A}_q \Phi_1^f & \cdots & \tilde{A}_q \Phi_q^f \end{pmatrix}, \tag{C5}
 \end{aligned}$$

597 and

$$[\tilde{X}; \tilde{X}] = \begin{pmatrix} z_q & \cdots & z_{n-2q} & z_{n-2q+1} & z_{n-2q+2} & \cdots & z_{n-q} & \boxed{z_1^f} & \cdots & z_{q-1}^f \\ z_{q+1} & \cdots & z_{n-2q+1} & z_{n-2q+2} & z_{n-2q+3} & \cdots & \boxed{z_1^f} & z_2^f & \cdots & z_q^f \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ z_{2q-2} & \cdots & z_{n-q} & z_{n-q+1} & z_{n-q} & \cdots & z_{q-1}^f & z_{q-1}^f & \cdots & z_{2q-2}^f \\ z_{2q-1} & \cdots & z_{n-q-1} & z_{n-q} & \boxed{z_1^f} & \cdots & z_q^f & z_q^f & \cdots & z_{2q-1}^f \end{pmatrix}. \tag{C6}$$

598 Comparing (C5) and (C6), it is clear that reconstructing the spatiotemporal patterns at  $s$  days in  
 599 the future requires only the prediction of the time series for  $s$  days.

## 600 References

601 Abhilash, S., A. Sahai, S. Pattnaik, B. Goswami, and A. Kumar, 2014a: Extended range prediction  
 602 of active-break spells of indian summer monsoon rainfall using an ensemble prediction system  
 603 in ncep climate forecast system. *International Journal of Climatology*, **34** (1), 98–113.

604 Abhilash, S., and Coauthors, 2014b: Prediction and monitoring of monsoon intraseasonal os-  
605 cillations over indian monsoon region in an ensemble prediction system using cfsv2. *Climate*  
606 *dynamics*, **42 (9-10)**, 2801–2815.

607 Belkin, M., and P. Niyogi, 2003: Laplacian eigenmaps for dimensionality reduction and data  
608 representation. *Neural computation*, **15 (6)**, 1373–1396.

609 Branicki, M., N. Chen, and A. J. Majda, 2013: Non-Gaussian test models for prediction and state  
610 estimation with model errors. *Chinese Annals of Mathematics, Series B*, **34 (1)**, 29–64.

611 Brenowitz, N. D., D. Giannakis, and A. J. Majda, 2016: Nonlinear laplacian spectral analysis of  
612 rayleigh–bénard convection. *Journal of Computational Physics*, **315**, 536–553.

613 Chen, N., and A. J. Majda, 2015a: Predicting the cloud patterns for the boreal summer intrasea-  
614 sonal oscillation through a low-order stochastic model. *Mathematics of Climate and Weather*  
615 *Forecasting*, **1 (1)**, 1–20.

616 Chen, N., and A. J. Majda, 2015b: Predicting the real-time multivariate madden–julian oscillation  
617 index through a low-order nonlinear stochastic model. *Monthly Weather Review*, **143 (6)**, 2148–  
618 2169.

619 Chen, N., and A. J. Majda, 2016: Filtering nonlinear turbulent dynamical systems through condi-  
620 tional gaussian statistics. *Monthly Weather Review*, **144 (12)**, 4885–4917.

621 Chen, N., A. J. Majda, and D. Giannakis, 2014: Predicting the cloud patterns of the madden-  
622 julian oscillation through a low-order nonlinear stochastic model. *Geophysical Research Letters*,  
623 **41 (15)**, 5612–5619.

624 Coifman, R. R., and S. Lafon, 2006: Diffusion maps. *Applied and computational harmonic anal-*  
625 *ysis*, **21 (1)**, 5–30.



626 Comeau, D., Z. Zhao, D. Giannakis, and A. J. Majda, 2016: Data-driven prediction strategies for  
627 low-frequency patterns of north pacific climate variability. *Climate Dynamics*, 1–18.

628 Crommelin, D., and A. Majda, 2004: Strategies for model reduction: comparing different optimal  
629 bases. *Journal of the Atmospheric Sciences*, **61** (17), 2206–2217.

630 Gadgil, S., 2003: The indian monsoon and its variability. *Annual Review of Earth and Planetary*  
631 *Sciences*, **31** (1), 429–467.

632 Giannakis, D., and A. J. Majda, 2011: Time series reconstruction via machine learning: Revealing  
633 decadal variability and intermittency in the north pacific sector of a coupled climate model.  
634 *CIDU*, Citeseer, 107–117.

635 Giannakis, D., and A. J. Majda, 2012a: Comparing low-frequency and intermittent variability  
636 in comprehensive climate models through nonlinear laplacian spectral analysis. *Geophysical*  
637 *Research Letters*, **39** (10).

638 Giannakis, D., and A. J. Majda, 2012b: Nonlinear laplacian spectral analysis for time series with  
639 intermittency and low-frequency variability. *Proceedings of the National Academy of Sciences*,  
640 **109** (7), 2222–2227.

641 Giannakis, D., and A. J. Majda, 2013: Nonlinear laplacian spectral analysis: capturing intermittent  
642 and low-frequency spatiotemporal patterns in high-dimensional data. *Statistical Analysis and*  
643 *Data Mining*, **6** (3), 180–194.

644 Giannakis, D., W.-w. Tung, and A. J. Majda, 2012: Hierarchical structure of the madden-julian  
645 oscillation in infrared brightness temperature revealed through nonlinear laplacian spectral anal-  
646 ysis. *Intelligent Data Understanding (CIDU), 2012 Conference on*, IEEE, 55–62.

647 Goswami, B., R. Ajayamohan, P. K. Xavier, and D. Sengupta, 2003: Clustering of synoptic activity  
648 by indian summer monsoon intraseasonal oscillations. *Geophysical Research Letters*, **30** (8).

649 Goswami, B., V. Krishnamurthy, and H. Annmalai, 1999: A broad-scale circulation index for the  
650 interannual variability of the indian summer monsoon. *Quarterly Journal of the Royal Meteorological Society*, **125** (554), 611–633.

651

652 Goswami, B., and R. A. Mohan, 2001: Intraseasonal oscillations and interannual variability of the  
653 indian summer monsoon. *Journal of Climate*, **14** (6), 1180–1198.

654 Harlim, J., A. Mahdi, and A. J. Majda, 2014: An ensemble Kalman filter for statistical estimation  
655 of physics constrained nonlinear regression models. *Journal of Computational Physics*, **257**,  
656 782–812.

657 Huffman, G. J., R. F. Adler, M. M. Morrissey, D. T. Bolvin, S. Curtis, R. Joyce, B. McGavock,  
658 and J. Susskind, 2001: Global precipitation at one-degree daily resolution from multisatellite  
659 observations. *Journal of Hydrometeorology*, **2** (1), 36–50.

660 Kikuchi, K., B. Wang, and Y. Kajikawa, 2012: Bimodal representation of the tropical intraseasonal  
661 oscillation. *Climate dynamics*, **38** (9-10), 1989–2000.

662 Kleeman, R., 2002: Measuring dynamical prediction utility using relative entropy. *Journal of the  
663 atmospheric sciences*, **59** (13), 2057–2072.

664 Kondrashov, D., M. Chekroun, A. Robertson, and M. Ghil, 2013: Low-order stochastic model and  
665 past-noise forecasting of the madden-julian oscillation. *Geophysical Research Letters*, **40** (19),  
666 5305–5310.

667 Kravtsov, S., D. Kondrashov, and M. Ghil, 2005: Multilevel regression modeling of nonlinear  
668 processes: Derivation and applications to climatic variability. *Journal of Climate*, **18 (21)**, 4404–  
669 4424.

670 Lau, W. K.-M., and D. E. Waliser, 2011: *Intraseasonal variability in the atmosphere-ocean climate*  
671 *system*. Springer Science & Business Media.

672 Lee, J.-Y., B. Wang, M. C. Wheeler, X. Fu, D. E. Waliser, and I.-S. Kang, 2013: Real-time multi-  
673 variate indices for the boreal summer intraseasonal oscillation over the asian summer monsoon  
674 region. *Climate Dynamics*, **40 (1-2)**, 493–509.

675 Liptser, R. S., and A. N. Shiryaev, 2001: *Statistics of Random Processes II: II. Applications*, Vol. 2.  
676 Springer.

677 Majda, A. J., and M. Branicki, 2012: Lessons in uncertainty quantification for turbulent dynamical  
678 systems. *Discrete Cont. Dyn. Systems*, **32 (9)**, 3133–3221.

679 Majda, A. J., and B. Gershgorin, 2010: Quantifying uncertainty in climate change science through  
680 empirical information theory. *Proceedings of the National Academy of Sciences*, **107 (34)**,  
681 14 958–14 963.

682 Majda, A. J., and B. Gershgorin, 2011: Improving model fidelity and sensitivity for complex sys-  
683 tems through empirical information theory. *Proceedings of the National Academy of Sciences*,  
684 **108 (25)**, 10 044–10 049.

685 Majda, A. J., and J. Harlim, 2013: Physics constrained nonlinear regression models for time series.  
686 *Nonlinearity*, **26 (1)**, 201–271.

687 Murakami, T., L.-X. Chen, and A. Xie, 1986: Relationship among seasonal cycles, low-frequency  
688 oscillations, and transient disturbances as revealed from outgoing longwave radiation data.  
689 *Monthly Weather Review*, **114 (8)**, 1456–1465.

690 Packard, N. H., J. P. Crutchfield, J. D. Farmer, and R. S. Shaw, 1980: Geometry from a time series.  
691 *Physical review letters*, **45 (9)**, 712.

692 Sabeerali, C., R. Ajayamohan, D. Giannakis, and A. J. Majda, 2017: Extraction and prediction  
693 of indices for monsoon intraseasonal oscillations: an approach based on nonlinear laplacian  
694 spectral analysis. *Climate Dynamics*, 1–20.

695 Sahai, A., and Coauthors, 2013: Simulation and extended range prediction of monsoon intrasea-  
696 sonal oscillations in ncep cfs/gfs version 2 framework. *Current Science*, **104 (10)**, 1394–1408.

697 Sauer, T., J. A. Yorke, and M. Casdagli, 1991: Embedology. *Journal of statistical Physics*, **65 (3)**,  
698 579–616.

699 Sikka, D., and S. Gadgil, 1980: On the maximum cloud zone and the itcz over indian, longitudes  
700 during the southwest monsoon. *Monthly Weather Review*, **108 (11)**, 1840–1853.

701 Slawinska, J., and D. Giannakis, 2016: Indo-pacific variability on seasonal to multidecadal  
702 timescales. part i: Intrinsic sst modes in models and observations. *arXiv preprint arX-*  
703 *iv:1604.01742*.

704 Suhas, E., J. Neena, and B. Goswami, 2013: An indian monsoon intraseasonal oscillations (miso)  
705 index for real time monitoring and forecast verification. *Climate dynamics*, **40 (11-12)**, 2605–  
706 2616.

- 707 Székely, E., D. Giannakis, and A. J. Majda, 2016a: Extraction and predictability of coherent  
708 intraseasonal signals in infrared brightness temperature data. *Climate Dynamics*, **46 (5-6)**,  
709 1473–1502.
- 710 Székely, E., D. Giannakis, and A. J. Majda, 2016b: Initiation and termination of intraseasonal  
711 oscillations in nonlinear laplacian spectral analysis-based indices. *Mathematics of Climate and*  
712 *Weather Forecasting*, **2 (1)**, 1–25.
- 713 Tung, W.-w., D. Giannakis, and A. J. Majda, 2014: Symmetric and antisymmetric convection  
714 signals in the madden–julian oscillation. part i: basic modes in infrared brightness temperature.  
715 *Journal of the Atmospheric Sciences*, **71 (9)**, 3302–3326.
- 716 Webster, P. J., V. O. Magana, T. Palmer, J. Shukla, R. Tomas, M. Yanai, and T. Yasunari, 1998:  
717 Monsoons: Processes, predictability, and the prospects for prediction. *Journal of Geophysical*  
718 *Research: Oceans*, **103 (C7)**, 14 451–14 510.
- 719 Wheeler, M. C., and H. H. Hendon, 2004: An all-season real-time multivariate mjo index: De-  
720 velopment of an index for monitoring and prediction. *Monthly Weather Review*, **132 (8)**, 1917–  
721 1932.
- 722 Zhang, C., 2005: Madden-julian oscillation. *Reviews of Geophysics*, **43 (2)**.

723 **LIST OF TABLES**

724 **Table 1.** Top: Optimal parameters for the nonlinear low-order stochastic model (1). The  
 725 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  
 726 units  $m^{-1/2}$ ;  $\phi$  is dimensionless. Here  $m$  stands for month. Bottom: the pa-  
 727 rameters from Chen and Majda (2015b), predicting the two boreal summer in-  
 728 traseasonal oscillation (BSISO) modes obtained by applying NLSA to the OLR  
 729 dataset, the Cloud Archive User Service (CLAUS) Version 4.7. The compari-  
 730 son between Section 6c. . . . . 39

	$d_u$	$d_v$	$d_\omega$	$\sigma_u$	$\sigma_v$	$\sigma_\omega$	$\gamma$	$a$	$f_0$	$f_I$	$\omega_f$	$\phi$
I). Optimal parameters	0.8	0.6	0.5	0.5	0.5	0.7	0.3	4.1	1.0	4.7	$2\pi/12$	-2
II). OLR-based parameters	0.9	0.9	0.5	0.3	0.8	1.0	0.3	4.25	1.0	4.0	$2\pi/12$	-1.4

731 TABLE 1. Top: Optimal parameters for the nonlinear low-order stochastic model (1). The parameters  
732  $d_u, a, \gamma, f_0, f_I, \omega_f, d_v$  and  $d_\omega$  have units  $m^{-1}$ ;  $\sigma_u, \sigma_v$  and  $\sigma_\omega$  have units  $m^{-1/2}$ ;  $\phi$  is dimensionless. Here  $m$  s-  
733 tands for month. Bottom: the parameters from Chen and Majda (2015b), predicting the two boreal summer  
734 intraseasonal oscillation (BSISO) modes obtained by applying NLSA to the OLR dataset, the Cloud Archive  
735 User Service (CLAUS) Version 4.7. The comparison between Section 6c.

**LIST OF FIGURES**

736

737 **Fig. 1.** Left: MISO indices from NLSA ranging from 1998/01/01 to 2013/12/31. The period up to 2007/12/31 is utilized as the training period in Section 3a for model calibration and the prediction skill is tested for year 2008 to 2013 in Section 4. Right: the highly non-Gaussian probability density functions (PDFs) with fat-tails of the MISO indices in the training period, where the sub-panel shows the PDF in logarithm scale. . . . . 42

738

739

740

741

742 **Fig. 2.** Comparison of the statistics of the MISO indices and those from the nonlinear low-order stochastic model (1)–(2) with optimal parameters in Table 1. (a) Long-term autorrelation and cross-correlation function up to 2 years. (b) Short-term autocorrelation function and cross-correlation functions up to 3 months. (c) PDFs of the MISO indices and the model signals. (d) Power spectrum of MISO indices and model signals. . . . . 43

743

744

745

746

747 **Fig. 3.** Skill scores with RMS error and pattern correlation for prediction utilizing ensemble mean in different years. The dashed line in the top panel shows the standard deviation of the equilibrium distribution associated with the MISO indices and that in the bottom panel shows the threshold of  $\text{Corr} = 0.5$ . Useful predictions are defined by 1) the RMS error in the prediction is less than the standard deviation of the truth at the equilibrium and 2) the pattern correlation between the predicted signal and the truth is above 0.5. . . . . 44

748

749

750

751

752

753 **Fig. 4.** Prediction of MISO 1 at 15- and 25-day lead utilizing ensemble mean. . . . . 45

754

755 **Fig. 5.** Long-range prediction of MISO 1 for 8 months starting from different dates. April 1 is a time at the transition between the quiescent phase and the active phase of the MISO indices; June 1 is a starting date in the active mature phase while October 1 is a starting date in the decaying phase of MISO activity. . . . . 46

756

757

758 **Fig. 6.** Phase diagrams of MISO 1 and MISO 2 (blue) and the corresponding predictions [ensemble mean (red) and 50 ensemble members (green)], starting from 2009/07/01, 2008/06/01 and 2013/06/01 and lasting for 30 days. Blue dots and small red rectangles indicate the truth and prediction for every 5 days. . . . . 47

759

760

761

762 **Fig. 7.** Skill score with pattern correlation (left) and RMS error (right) for predicting the reconstructed spatiotemporal patterns as a function of lead time (days) for July 2009 (black), June 2008 (red) and June 2013 (green). . . . . 48

763

764

765 **Fig. 8.** Reconstruction and prediction of the spatiotemporal patterns of July 2009 starting from 1 July 2009. The prediction at day 1, 4, 8, 12, 16, 20, 24 and 30 are shown. . . . . 49

766

767 **Fig. 9.** Same as Figure 8 but for June 2008. . . . . 50

768 **Fig. 10.** Comparison of the truth (first row), the approximated spatiotemporal patterns from NLSA based on (9) (second row) and the patterns in which the spatial basis is obtained based on the modes by applying running average to the raw data as described in Section 5 (third row). . . . . 51

769

770

771 **Fig. 11.** Comparison of the statistics based on the full period of MISO indices (1998-2013; black), training period utilized in Section 3a (1998-2007; red) and the 3-year short training period in Section 6a (1998-2000; green). . . . . 52

772

773

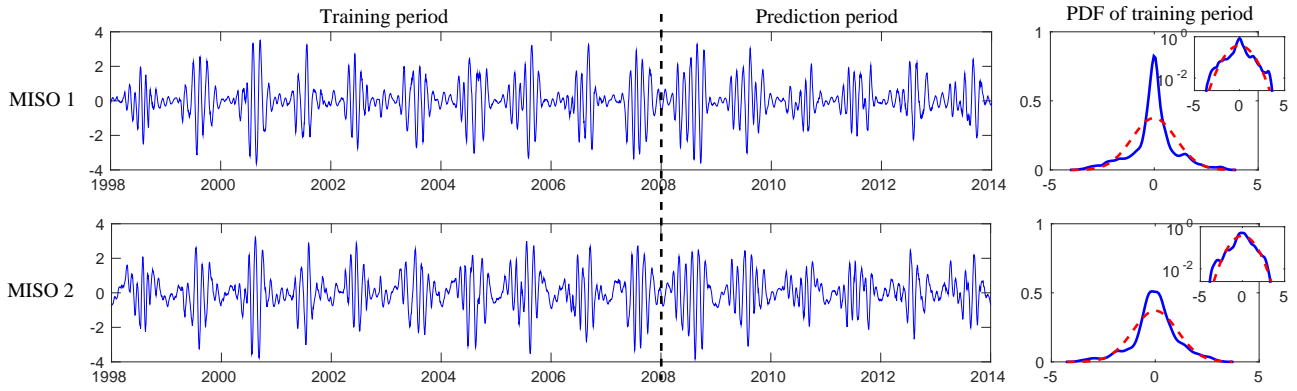
774 **Fig. 12.** Left: Skill score with RMS error and pattern correlation for predicting the MISO indices in different years from 2001 to 2007. As in Figure 3, the two dashes lines indicate the standard deviation of the MISO indices at climatology and the value with  $\text{Corr} = 0.5$ , which serve

775

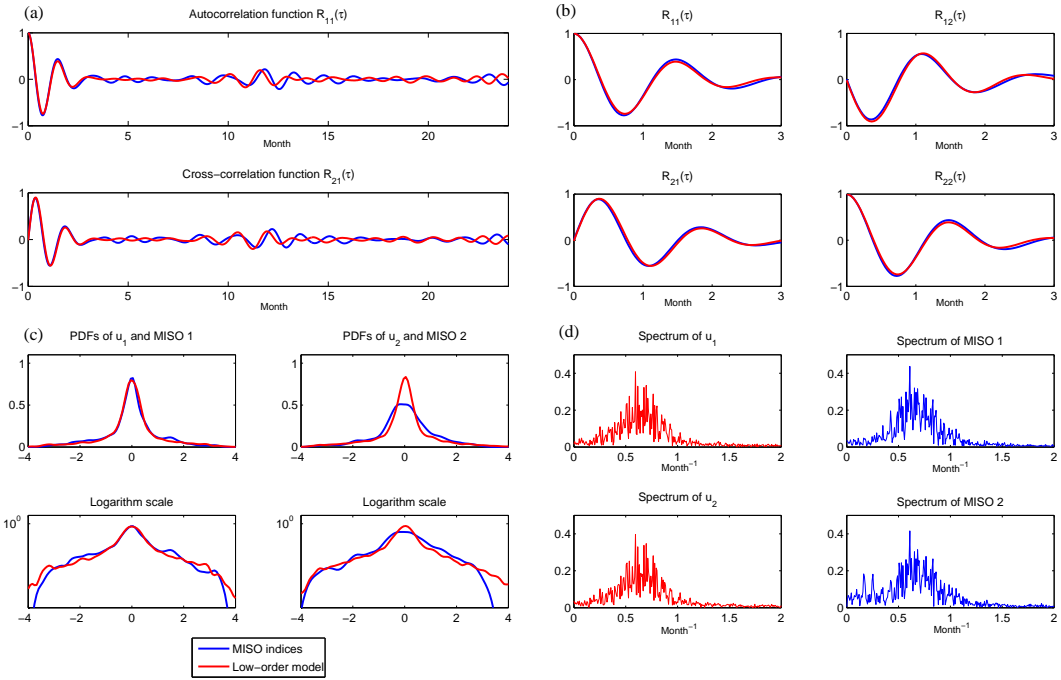
776



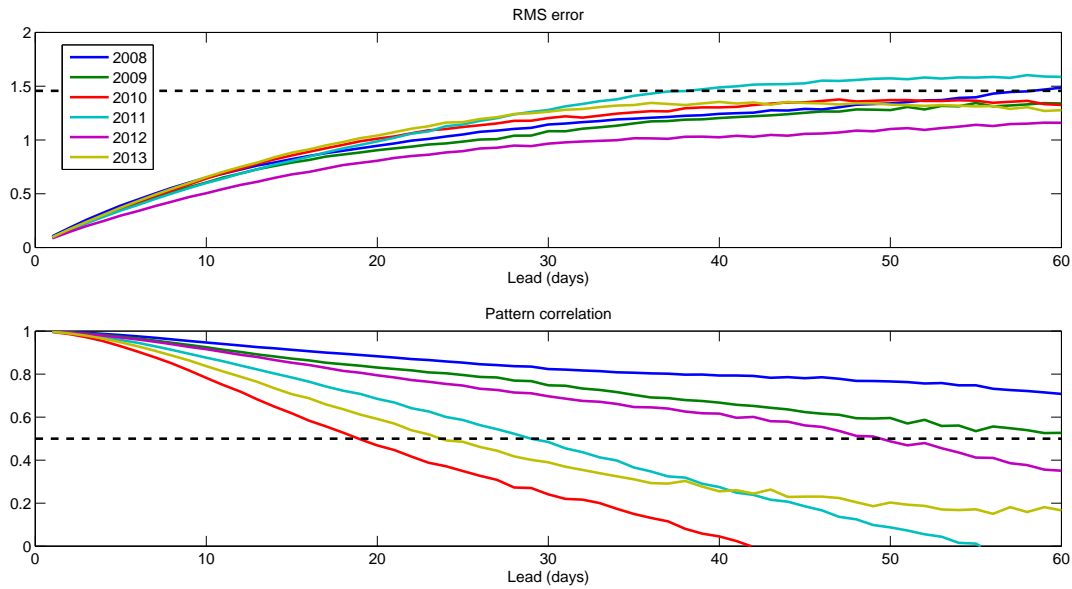
777	as the threshold for the useful prediction. Right: 25-day lead prediction of MISO 1 at four	
778	different years. The model parameters are listed in the top row of Table 1. . . . .	53
779	<b>Fig. 13.</b> Comparison of the indices, obtained by applying NLSA with different lagged embedding	
780	window sizes ( $q = 64, 48$ and $34$ days), and the associated statistics. . . . .	54
781	<b>Fig. 14.</b> Comparison of the prediction skill of the indices obtained by applying NLSA with different	
782	lagged embedding window sizes (top: $q = 64$ ; middle: $q = 48$ ; bottom: $q = 34$ ). Left: useful	
783	prediction for the full year and for only the boreal summer time (June to September). Right:	
784	15-day lead prediction for the indices obtained with different lagged embedding window	
785	sizes. . . . .	55
786	<b>Fig. 15.</b> Applying the parameters in Chen and Majda (2015b) to calibrate and predict the MISO	
787	precipitation indices. Panel (a) is the autocorrelation of $u_1$ up to 2 years. Panels (b) are for	
788	the autocorrelation of $u_1$ and crosscorrelation between $u_1$ and $u_2$ up to 3 months. Panels (c)	
789	are for the PDF in linear and logarithm scales, respectively. Panel (d) compares the power	
790	spectrum of $u_1$ . Panels (e) and (f) compare the 25-day lead ensemble mean predictions with	
791	the optimal parameters and the parameters borrowed from OLR dataset in Chen and Majda	
792	(2015b). . . . .	56
793	<b>Fig. A1.</b> Sensitive test. Left: Model error (via information distance) as functions of different param-	
794	eters. Right: RMS error in the autocorrelation functions as functions of those parameters	
795	related to the phase. The black dots indicate the optimal parameters as listed in Table 1. The	
796	two dotted lines in each panel indicate the range of randomly-picked suboptimal parameters	
797	in the test of prediction. . . . .	57



798 FIG. 1. Left: MISO indices from NLSA ranging from 1998/01/01 to 2013/12/31. The period up to 2007/12/31  
 799 is utilized as the training period in Section 3a for model calibration and the prediction skill is tested for year 2008  
 800 to 2013 in Section 4. Right: the highly non-Gaussian probability density functions (PDFs) with fat-tails of the  
 801 MISO indices in the training period, where the sub-panel shows the PDF in logarithm scale.



802 FIG. 2. Comparison of the statistics of the MISO indices and those from the nonlinear low-order stochastic  
 803 model (1)–(2) with optimal parameters in Table 1. (a) Long-term autorrelation and cross-correlation function up  
 804 to 2 years. (b) Short-term autocorrelation function and cross-correlation functions up to 3 months. (c) PDFs of  
 805 the MISO indices and the model signals. (d) Power spectrum of MISO indices and model signals.



806 FIG. 3. Skill scores with RMS error and pattern correlation for prediction utilizing ensemble mean in different  
 807 years. The dashed line in the top panel shows the standard deviation of the equilibrium distribution associated  
 808 with the MISO indices and that in the bottom panel shows the threshold of  $\text{Corr} = 0.5$ . Useful predictions are  
 809 defined by 1) the RMS error in the prediction is less than the standard deviation of the truth at the equilibrium  
 810 and 2) the pattern correlation between the predicted signal and the truth is above 0.5.

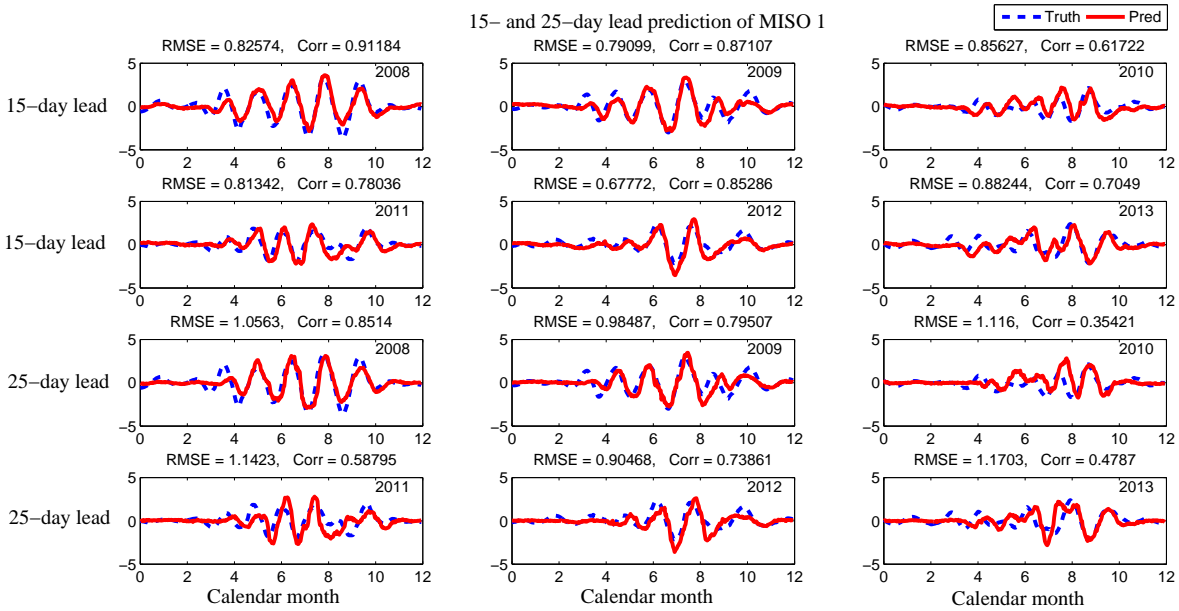
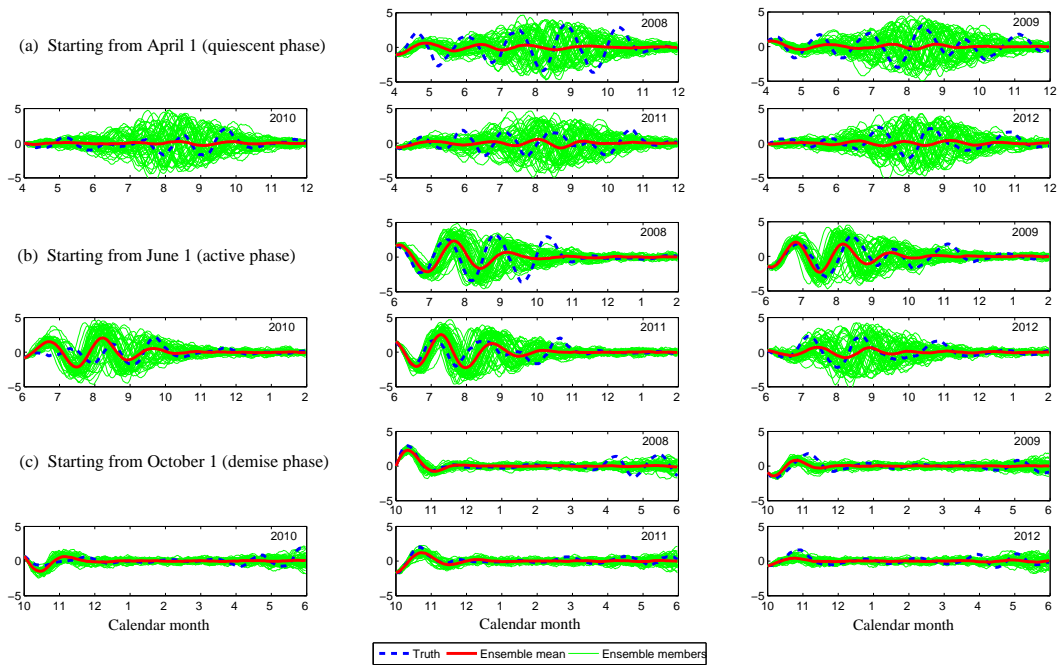
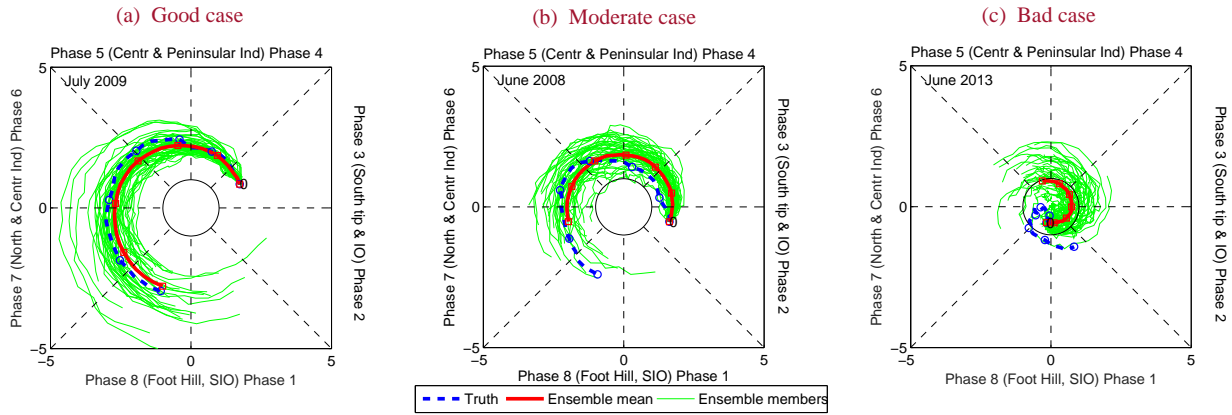


FIG. 4. Prediction of MISO 1 at 15- and 25-day lead utilizing ensemble mean.

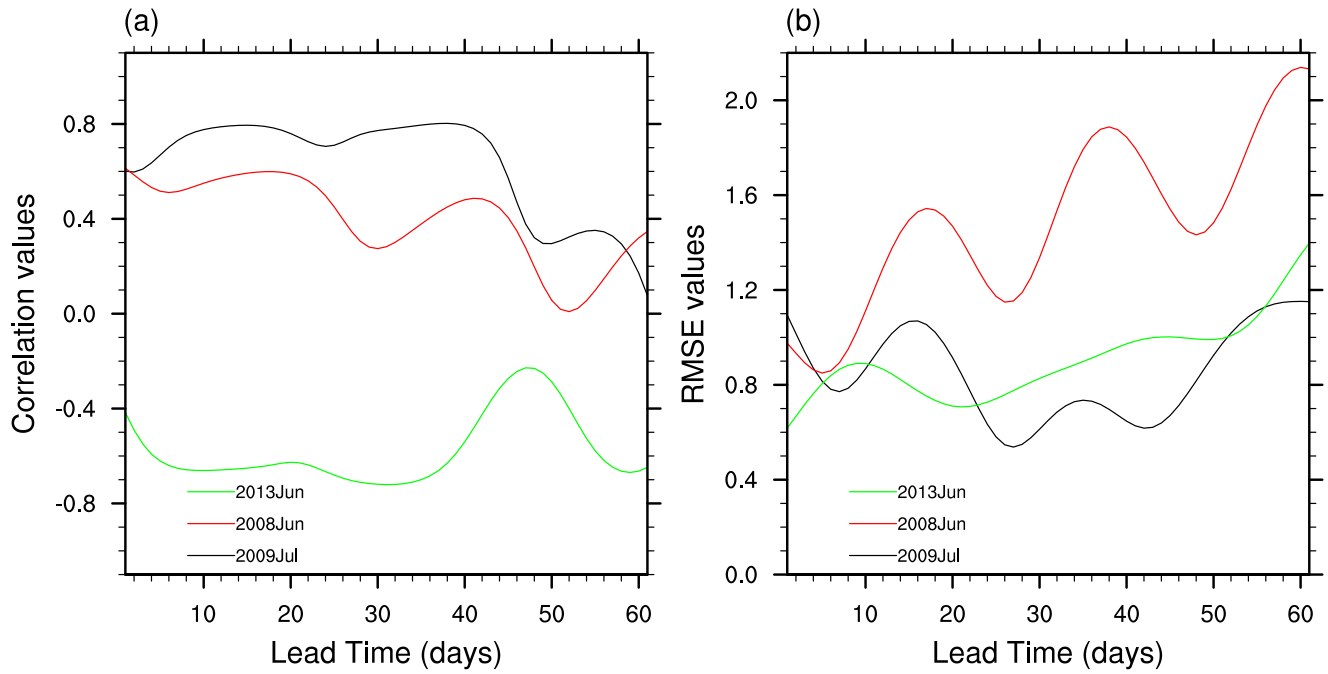
Medium- and long-range forecasting



811 FIG. 5. Long-range prediction of MISO 1 for 8 months starting from different dates. April 1 is a time at the  
 812 transition between the quiescent phase and the active phase of the MISO indices; June 1 is a starting date in the  
 813 active mature phase while October 1 is a starting date in the decaying phase of MISO activity.

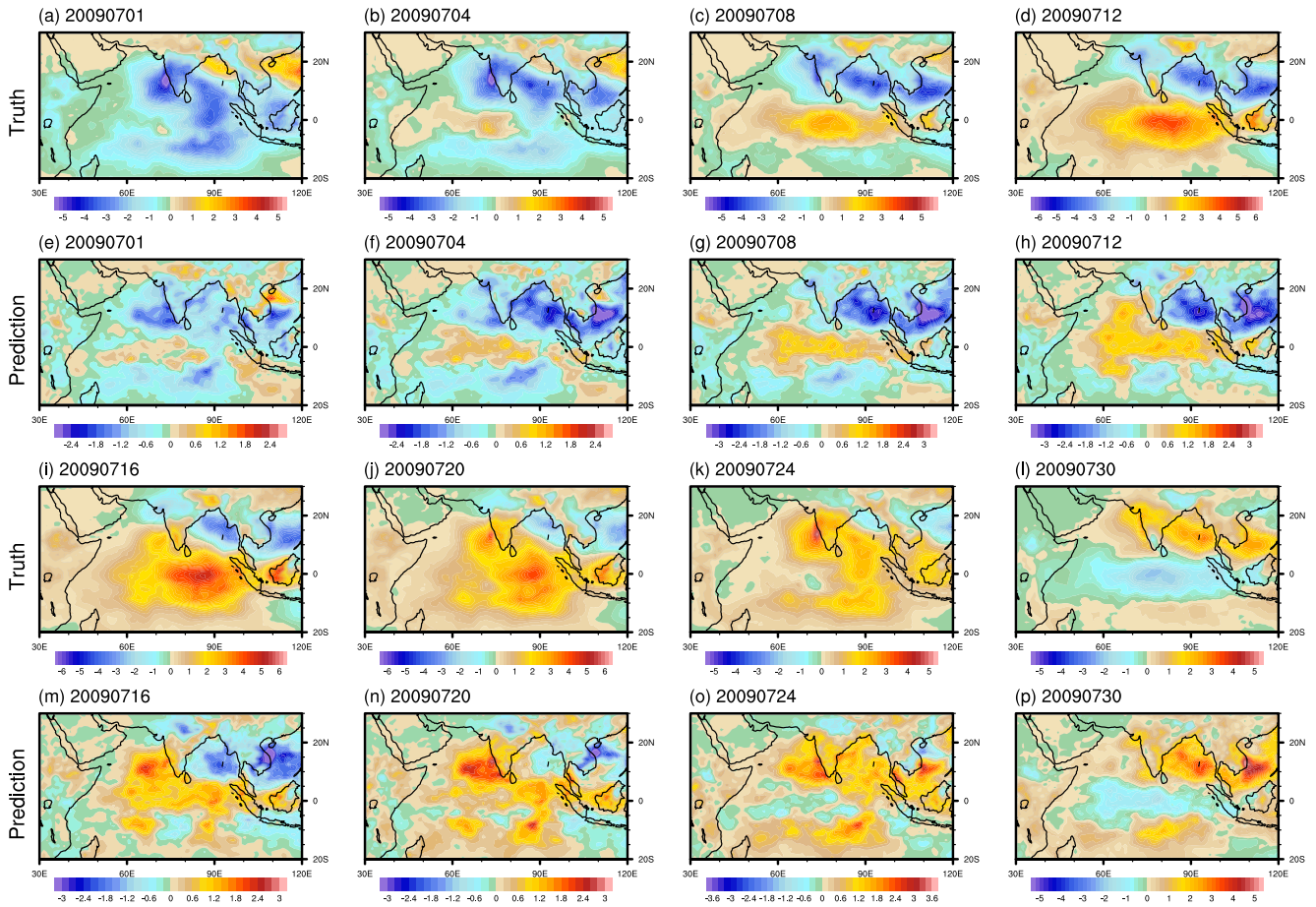


814 FIG. 6. Phase diagrams of MISO 1 and MISO 2 (blue) and the corresponding predictions [ensemble mean  
 815 (red) and 50 ensemble members (green)], starting from 2009/07/01, 2008/06/01 and 2013/06/01 and lasting for  
 816 30 days. Blue dots and small red rectangles indicate the truth and prediction for every 5 days.



817 FIG. 7. Skill score with pattern correlation (left) and RMS error (right) for predicting the reconstructed  
 818 spatiotemporal patterns as a function of lead time (days) for July 2009 (black), June 2008 (red) and June 2013  
 819 (green).





820 FIG. 8. Reconstruction and prediction of the spatiotemporal patterns of July 2009 starting from 1 July 2009.  
 821 The prediction at day 1, 4, 8, 12, 16, 20, 24 and 30 are shown.

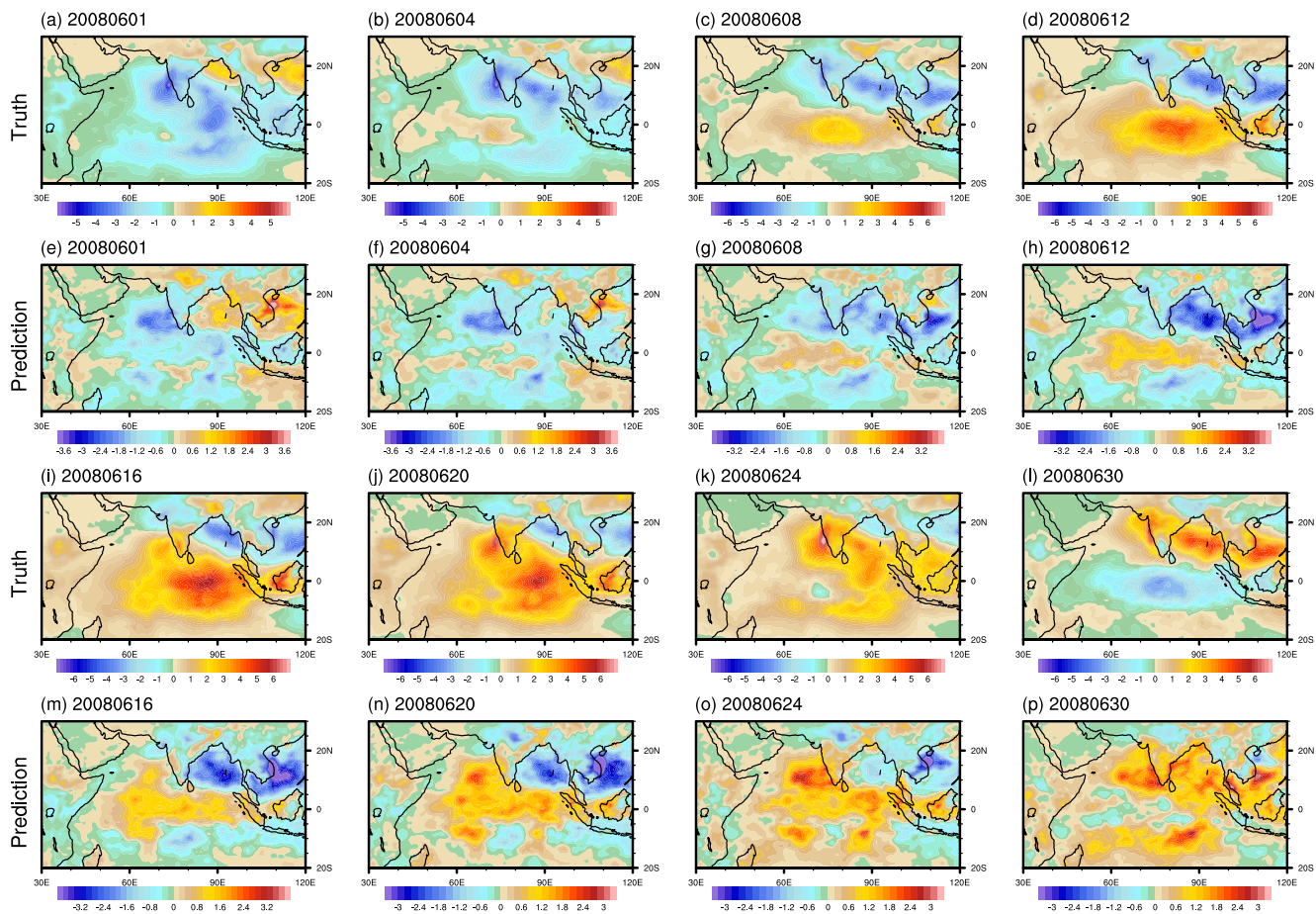
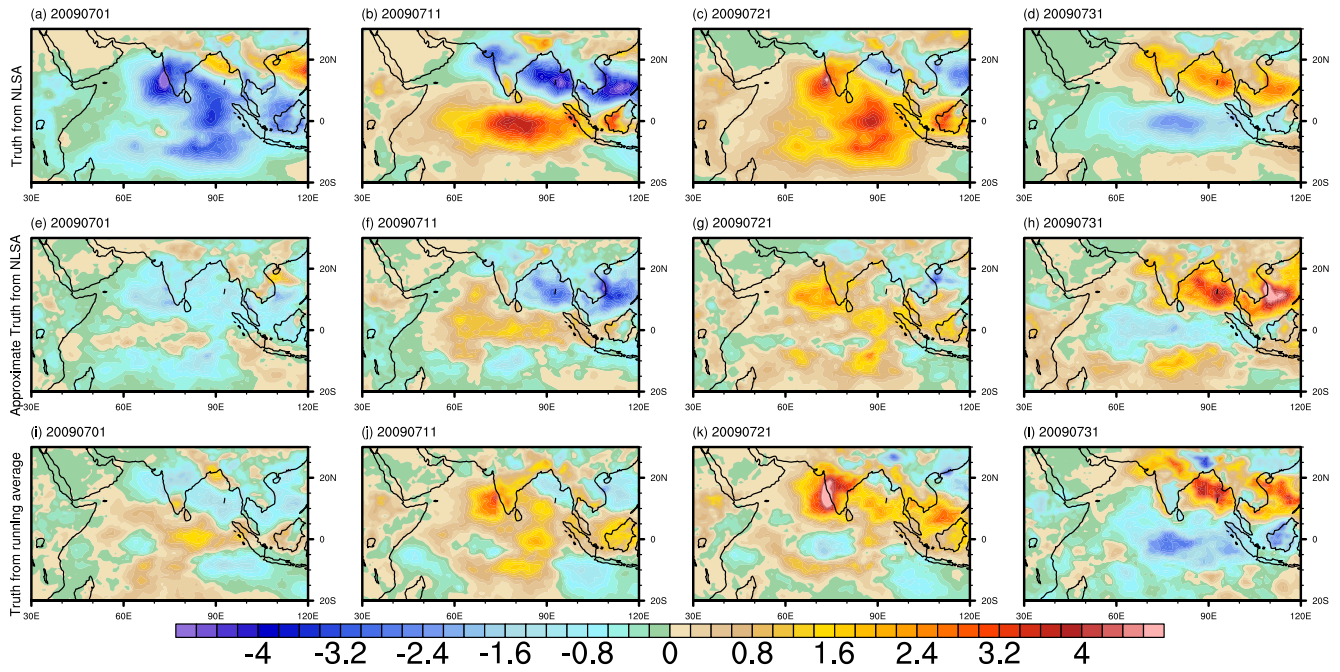
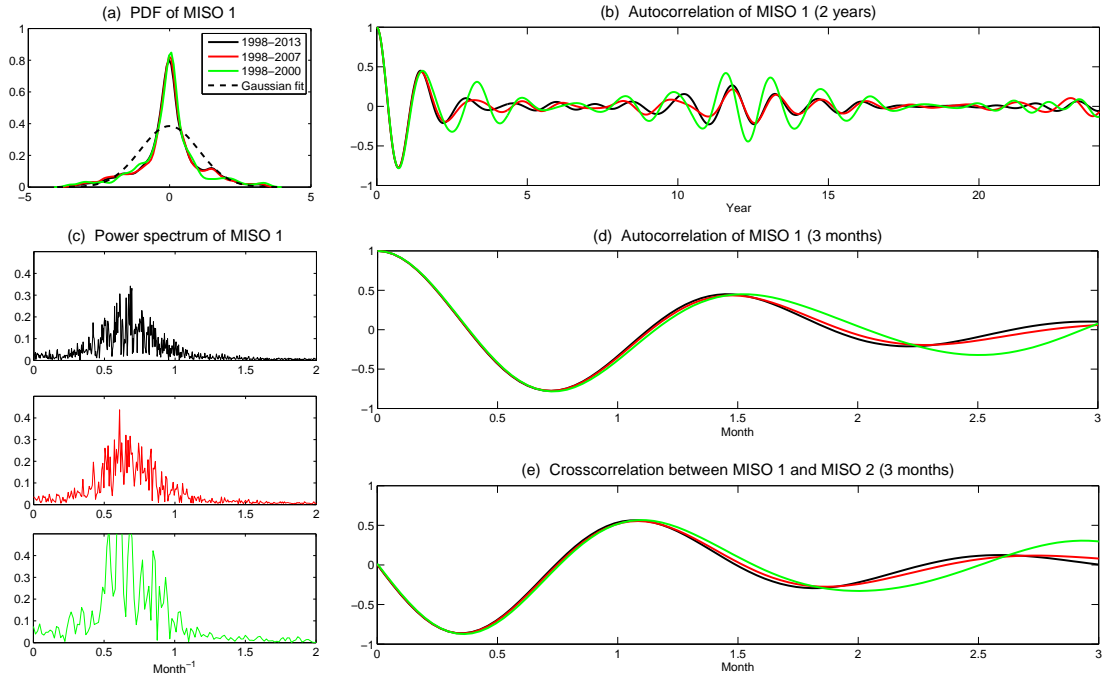


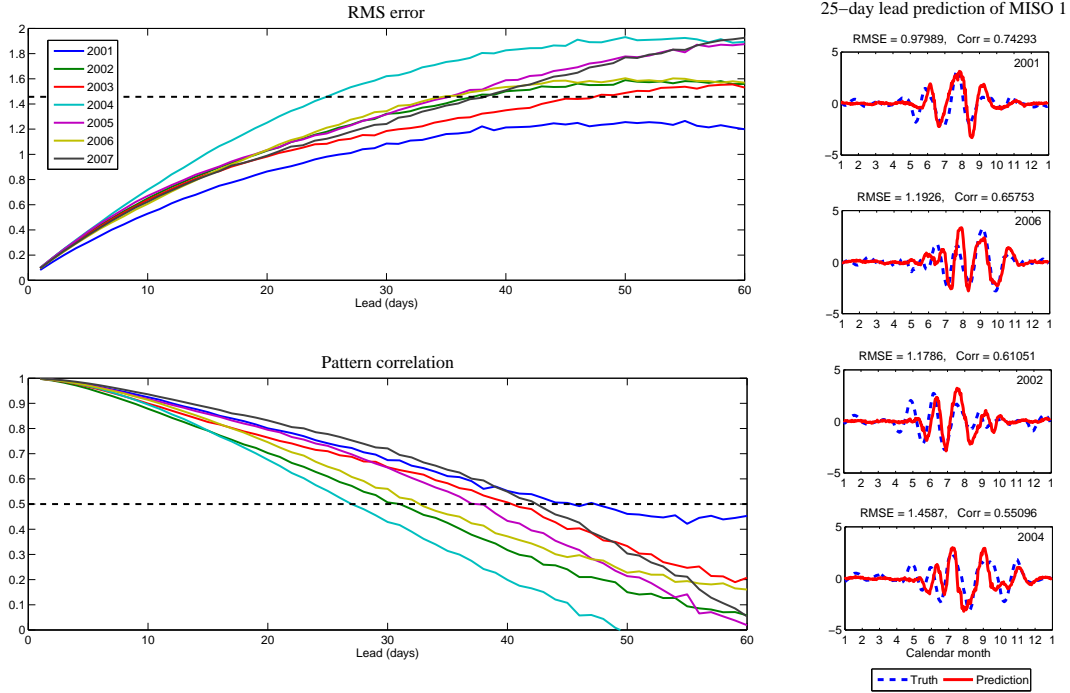
FIG. 9. Same as Figure 8 but for June 2008.



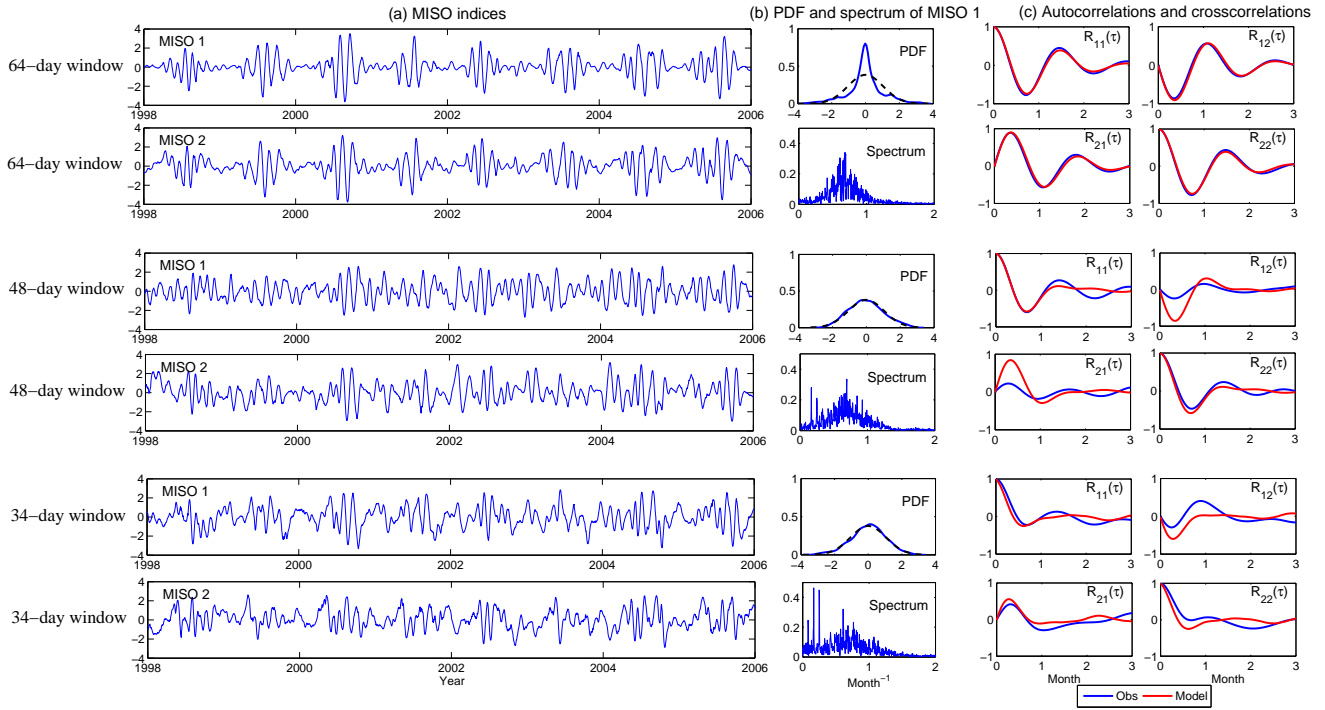
822 FIG. 10. Comparison of the truth (first row), the approximated spatiotemporal patterns from NLSA based on  
 823 (9) (second row) and the patterns in which the spatial basis is obtained based on the modes by applying running  
 824 average to the raw data as described in Section 5 (third row).



825 FIG. 11. Comparison of the statistics based on the full period of MISO indices (1998-2013; black), training  
 826 period utilized in Section 3a (1998-2007; red) and the 3-year short training period in Section 6a (1998-2000;  
 827 green).

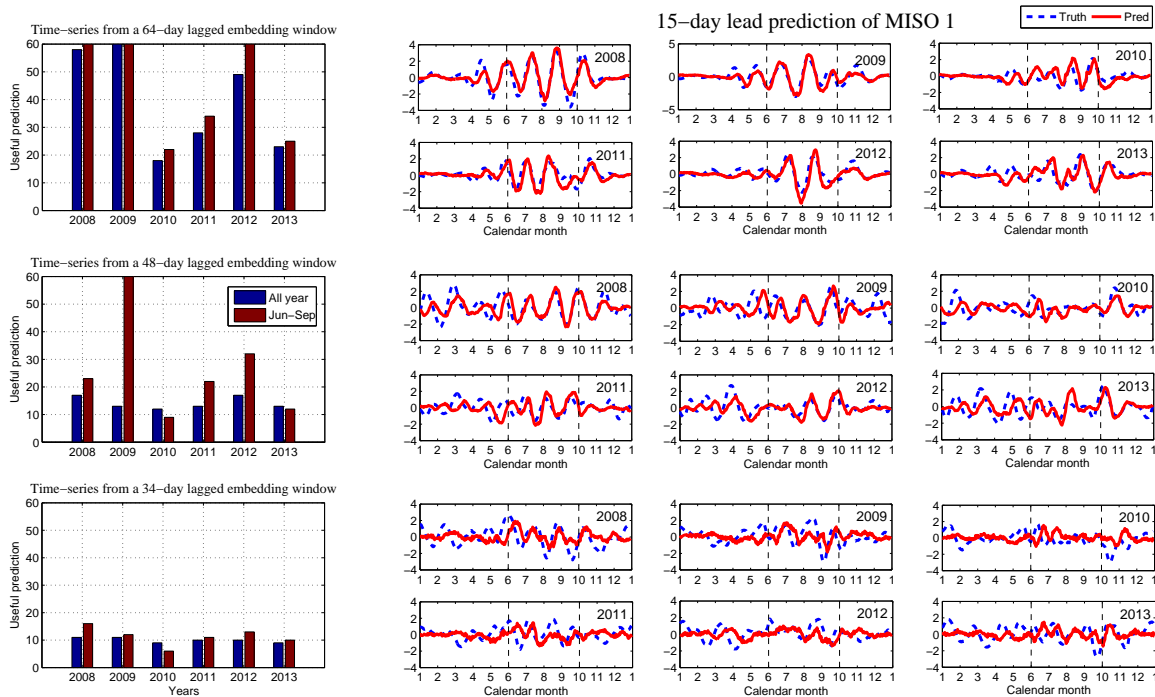


828 FIG. 12. Left: Skill score with RMS error and pattern correlation for predicting the MISO indices in different  
 829 years from 2001 to 2007. As in Figure 3, the two dashes lines indicate the standard deviation of the MISO  
 830 indices at climatology and the value with  $\text{Corr} = 0.5$ , which serve as the threshold for the useful prediction.  
 831 Right: 25-day lead prediction of MISO 1 at four different years. The model parameters are listed in the top row  
 832 of Table 1.

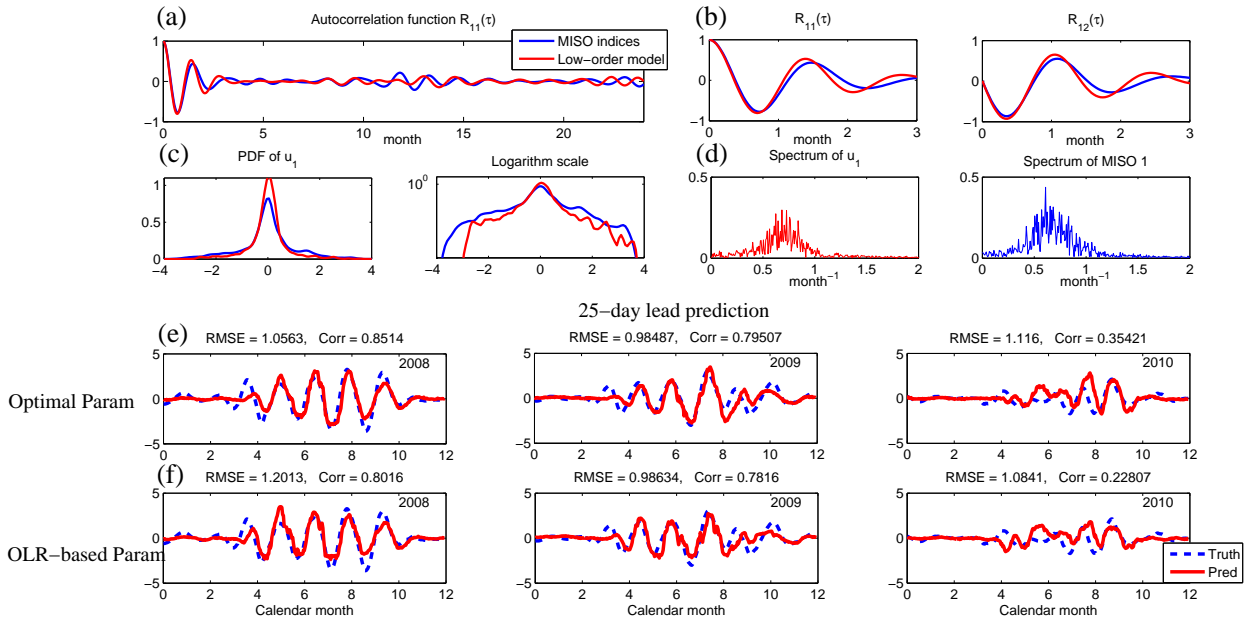


833 FIG. 13. Comparison of the indices, obtained by applying NLSA with different lagged embedding window  
 834 sizes ( $q = 64, 48$  and  $34$  days), and the associated statistics.



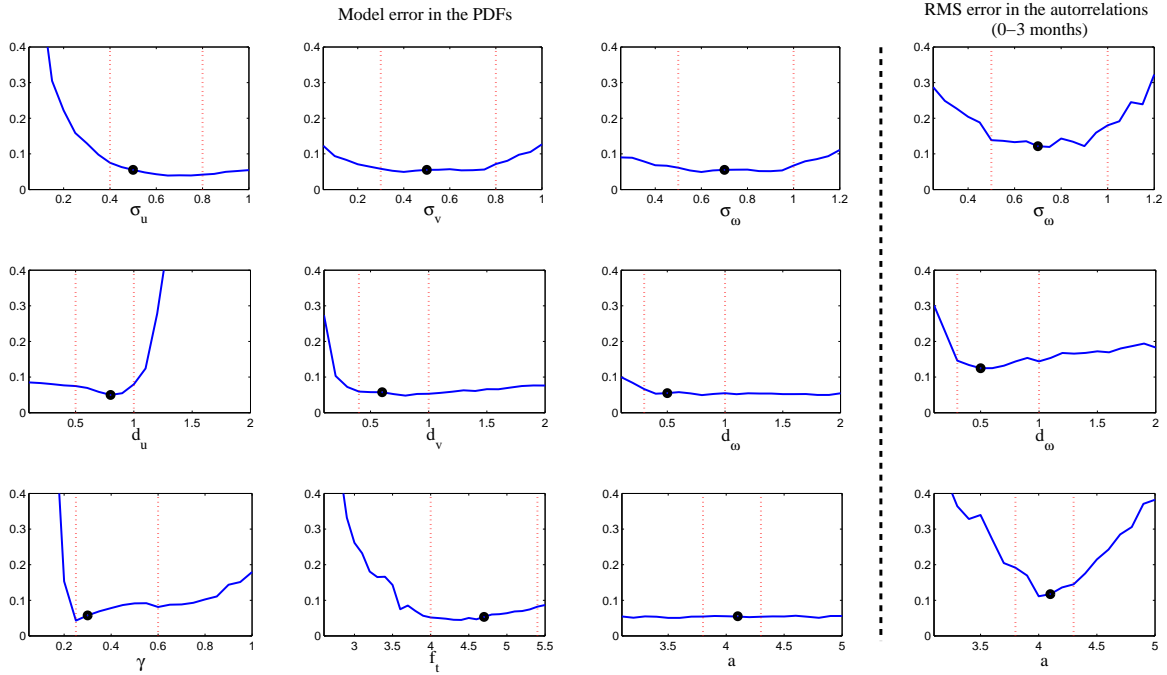


835 FIG. 14. Comparison of the prediction skill of the indices obtained by applying NLSA with different lagged  
 836 embedding window sizes (top:  $q = 64$ ; middle:  $q = 48$ ; bottom:  $q = 34$ ). Left: useful prediction for the full year  
 837 and for only the boreal summer time (June to September). Right: 15-day lead prediction for the indices obtained  
 838 with different lagged embedding window sizes.



839 FIG. 15. Applying the parameters in Chen and Majda (2015b) to calibrate and predict the MISO precipitation  
 840 indices. Panel (a) is the autocorrelation of  $u_1$  up to 2 years. Panels (b) are for the autocorrelation of  $u_1$  and  
 841 crosscorrelation between  $u_1$  and  $u_2$  up to 3 months. Panels (c) are for the PDF in linear and logarithm scales,  
 842 respectively. Panel (d) compares the power spectrum of  $u_1$ . Panels (e) and (f) compare the 25-day lead ensemble  
 843 mean predictions with the optimal parameters and the parameters borrowed from OLR dataset in Chen and  
 844 Majda (2015b).





845 Fig. A1. Sensitive test. Left: Model error (via information distance) as functions of different parameters.  
 846 Right: RMS error in the autocorrelation functions as functions of those parameters related to the phase. The  
 847 black dots indicate the optimal parameters as listed in Table 1. The two dotted lines in each panel indicate the  
 848 range of randomly-picked suboptimal parameters in the test of prediction.