

# Geophysical Research Letters

## RESEARCH LETTER

10.1029/2018GL081100

### Key Points:

- Singular spectrum analysis (SSA) and extended empirical orthogonal function (EEOF) methods suffer from endpoint issues
- SSA with conditional predictions (SSA-CP) is presented as a simple modification to improve real-time estimates near endpoints
- Forecasts are also possible, including error estimates, and are optimal for Gaussian data and shown to also be skillful for non-Gaussian data

### Supporting Information:

- Supporting Information S1

### Correspondence to:

H. Reed Ogrosky,  
hrogrosky@vcu.edu

### Citation:

Ogrosky, H. R., Stechmann, S. N., Chen, N., & Majda, A. J. (2019). Singular spectrum analysis with conditional predictions for real-time state estimation and forecasting. *Geophysical Research Letters*, 46. <https://doi.org/10.1029/2018GL081100>

Received 26 OCT 2018

Accepted 30 JAN 2019

Accepted article online 4 FEB 2019

## Singular Spectrum Analysis With Conditional Predictions for Real-Time State Estimation and Forecasting

H. Reed Ogrosky<sup>1</sup> , Samuel N. Stechmann<sup>2,3</sup> , Nan Chen<sup>2</sup>, and Andrew J. Majda<sup>4,5</sup>

<sup>1</sup>Department of Mathematics and Applied Mathematics, Virginia Commonwealth University, Richmond, VA, USA, <sup>2</sup>Department of Mathematics, University of Wisconsin-Madison, Madison, WI, USA, <sup>3</sup>Department of Atmospheric and Oceanic Sciences, University of Wisconsin-Madison, Madison, WI, USA, <sup>4</sup>Department of Mathematics and Center for Atmosphere Ocean Science, Courant Institute of Mathematical Sciences, New York University, New York, NY, USA, <sup>5</sup>Center for Prototype Climate Modeling, NYU Abu Dhabi, Abu Dhabi, United Arab Emirates

**Abstract** Singular spectrum analysis (SSA) or extended empirical orthogonal function methods are powerful, commonly used data-driven techniques to identify modes of variability in time series and space-time data sets. Due to the time-lagged embedding, these methods can provide inaccurate reconstructions of leading modes near the endpoints, which can hinder the use of these methods in real time. A modified version of the traditional SSA algorithm, referred to as SSA with conditional predictions (SSA-CP), is presented to address these issues. It is tested on low-dimensional, approximately Gaussian data, high-dimensional non-Gaussian data, and partially observed data from a multiscale model. In each case, SSA-CP provides a more accurate real-time estimate of the leading modes of variability than the traditional reconstruction. SSA-CP also provides predictions of the leading modes and is easy to implement. SSA-CP is optimal in the case of Gaussian data, and the uncertainty in real-time estimates of leading modes is easily quantified.

**Plain Language Summary** Singular spectrum analysis (SSA) is a powerful, commonly used technique to identify prominent patterns in observed data. However, SSA has some difficulty in providing accurate estimates near the endpoints of the time series, which can hinder its use in real time. A modified version of the SSA algorithm, referred to as SSA with conditional predictions, is presented to address these issues. SSA with conditional predictions provides a more accurate real-time estimate of the leading modes of variability than the traditional method in a variety of tests. It can also be used to predict these patterns, and it is easy to implement. The uncertainty in the real-time estimates of leading patterns is easily quantified as well.

### 1. Introduction

Singular spectrum analysis (SSA) or extended empirical orthogonal function (EEOF) methods are powerful, commonly used tools available for identifying modes of variability in time series and space-time data sets. SSA's usefulness has been demonstrated in a variety of fields over the last 3–4 decades, including, for example, nonlinear dynamics (e.g., Broomhead & King, 1986), geoscience (e.g., Keppenne & Ghil, 1990; Kikuchi & Wang, 2008; Mo, 2001; Roundy & Schreck III, 2009; Weare & Nasstrom, 1982; Vautard & Ghil, 1989; Vautard et al., 1992), and economics (e.g., Hassani et al., 2014; Lisi & Medio, 1997). Its popularity is due both to its ease of implementation and to its ability to eliminate noise and extract trends, oscillations, and other signals in both univariate and multivariate time series.

As with some other methods for mode identification in space-time data (e.g., Fourier filtering), SSA suffers from endpoint issues; that is, estimates of leading modes can be inaccurate in real-time without future information. Therefore, SSA may provide inaccurate initial conditions for real-time forecasts. Despite these challenges, it is sometimes used either as a filtering step prior to generating real-time forecasts (e.g., Hassani et al., 2014; Mo, 2001) or in tests of forecast models (e.g., Chen & Majda, 2016; Kang & Kim, 2010; Kondrashov et al., 2013), due to its effectiveness at mode identification.

This motivates the question: Is there a modified version of SSA that (i) is as straightforward to implement as SSA but that (ii) provides the most accurate real-time state estimation possible of leading modes of variability?

This question, along with the related question of how to best modify SSA for use on data sets with gaps in the data, has motivated the proposal and study of numerous modified versions of SSA. These methods include schemes for modifying incomplete columns of the lag-embedded matrix by weighting known values (Schoellhamer, 2001), iterative SSA methods (Kondrashov & Ghil, 2006; Kondrashov et al., 2010), methods based on linear recurrent formulas (Golyandina & Osipov, 2007), methods that project smoothed data onto leading SSA modes computed with Fourier-filtered data (Roundy & Schreck III, 2009), combined recurrent forecasting and hindcasting (Rodrigues & de Carvalho, 2013), energy-minimizing reconstructions of principal components (PCs; Shen et al., 2014, 2015), and a method utilizing a predicted spatial basis (Chen et al., 2018). Some of these methods will be discussed in section 5.

Here we propose and study yet another modification of SSA. This method makes use of conditional mean predictions based on the covariance matrix of the lag-embedded data, and we refer to it as SSA with conditional predictions (SSA-CP). Another appropriate name would be real-time SSA, though we use SSA-CP here to avoid confusion with other methods proposed for using SSA in real time, some of which are discussed further in section 5.

The results of tests shown here suggest that this method is effective at addressing these endpoint issues in a variety of settings. The data sets used in these tests include both univariate data sets and multivariate data sets with small (2–3) or somewhat large (64) number of spatial dimensions, partially observed systems and data sets with all dynamical variables observed, Gaussian and non-Gaussian data, and synthetic time series and time series generated by observational data.

Given these results, there are at least four reasons for using this method. First, it is simple and easy to implement, requiring only small additional steps during the normal SSA algorithm. Second, it provides both state estimation and prediction of leading modes of variability. Third, it provides an optimal reconstruction if the data are Gaussian using the statistics of the first two moments. Fourth, it outperforms many other proposed methods of SSA state estimation for both Gaussian and non-Gaussian data.

The rest of the paper is organized as follows: Section 2 describes the traditional SSA method and the proposed modification. Section 3 lists data sets and models used in tests of this method. Results are presented in section 4. Discussion of the methods and results is given in section 5, including a brief comparison of the results with those of other modified SSA methods. Conclusions are given in section 6.

## 2. SSA Algorithms

A brief review of the traditional SSA algorithm is now given, followed by a description of the proposed modification. When used on multivariate time series, SSA is often referred to as multichannel SSA (MSSA) in the literature; here SSA will be used to refer to either the univariate or multivariate cases. The theory of SSA, which has been developed over the last several decades, is not discussed here; see, for example, (Aubry et al., 1991; Ghil et al., 2002; Golyandina et al., 2001; Hassani, 2007) for discussion of this underlying theory.

### 2.1. Traditional SSA

We briefly describe the traditional SSA algorithm for a data set with spatial dimension  $D$ ; the traditional univariate SSA algorithm can be reproduced by setting  $D = 1$  below.

Let  $\vec{x}_i$  be a  $D$ -dimensional column vector at time  $i$ , with  $1 \leq i \leq N$ . The four steps of SSA are as follows:

Step 1. Create the time-lagged embedding matrix  $\mathbf{X}$  of size  $(MD) \times (N - M + 1)$ :

$$\mathbf{X} = \begin{bmatrix} \vec{x}_1 & \vec{x}_2 & \dots & \vec{x}_{N-M+1} \\ \vec{x}_2 & \vec{x}_3 & \dots & \vec{x}_{N-M+2} \\ \vdots & \vdots & & \vdots \\ \vec{x}_{M-1} & \vec{x}_M & \dots & \vec{x}_{N-1} \\ \vec{x}_M & \vec{x}_{M+1} & \dots & \vec{x}_N \end{bmatrix}, \quad (1)$$

where  $M$  is the length of the embedding window.

Step 2. Find eigenvalues and eigenvectors of the covariance matrix  $\mathbf{C} = \mathbf{X}\mathbf{X}^T / (N - M + 1)$ . Each eigenvector  $\vec{v}$  (sometimes referred to as an EOF) is an  $(MD)$ -dimensional column vector with corresponding eigenvalue  $\lambda$ :

$$\vec{v} = [\vec{v}_1^T, \dots, \vec{v}_M^T]^T, \quad (2)$$

where  $\vec{v}_s$  is a  $D$ -dimensional column vector used to denote the lag- $s$  portion of the eigenvector.

Step 3. Find the PC of each mode by projecting the lag-embedded data onto the appropriate eigenvector:

$$\vec{\phi} = \mathbf{X}^T \vec{v}. \quad (3)$$

The entries of each PC will be denoted  $\vec{\phi} = [\phi_1, \dots, \phi_{N-M+1}]^T$ .

Step 4. Reconstruct the data corresponding to each mode by calculating the reconstructed component (RC)  $\vec{z}(t)$ :

$$\vec{z}(t) = \frac{1}{M_t} \sum_{i=L_t}^{U_t} \phi_{t-i+1} \vec{v}_i, \quad (4)$$

where  $(M_t, L_t, U_t)$  are defined by (see, e.g., Ghil et al., 2002)

$$(M_t, L_t, U_t) = \begin{cases} (t, 1, t), & 1 \leq t \leq M-1, \\ (M, 1, M), & M \leq t \leq N-M+1, \\ (N-t+1, t-N+M, M), & N-M+2 \leq t \leq N, \end{cases} \quad (5)$$

so that each RC  $\vec{z}$  is a (possibly multivariate) time series of length  $N$ , with each  $\vec{z}(t)$  a  $D$ -dimensional column vector.

Each RC entry at time  $t^*$  depends directly on one embedding window of PC entries, and each PC entry depends on one embedding window of data. As a result, each RC entry at time  $t^*$  is influenced primarily by the values of  $\vec{x}_{t^*-M+1}$  through  $\vec{x}_{t^*+M-1}$ ; that is, two embedding windows worth of data, spanning the window immediately prior to  $t^*$  and the window immediately following  $t^*$ , contribute directly to the reconstruction at  $t^*$ . For  $t^* > N-M$ , the embedding window's worth of data immediately following  $t^*$  is not entirely known. The reconstruction process makes use of the known data by averaging over the available products  $\phi_{t-i+1} \vec{v}_i$  in (4), but these final  $M-1$  entries of each reconstruction are only estimates of the state of each mode, and can be expected to change as data becomes available at times occurring after the end of the time series. (The same endpoint issues affect the reconstruction for  $t^* < M$ .)

The reconstruction method in (4) has been shown to be an optimal method, in the sense that, for  $D=1$ , for example, it produces the Hankel matrix that is closest to the matrix  $\vec{\phi} \vec{v}^T$  in matrix norm (Golyandina et al., 2001; Hassani, 2007). However, other reconstruction formulas may be considered, including ones that avoid the endpoint issues of the traditional reconstruction in (4). One such method is the “predicted spatial basis” method of Chen et al. (2018), in which a method that shifts future information to the spatial basis (and not the PCs) is tested on a monsoon intraseasonal oscillation index. The method in (4) is used as the primary basis for comparison here due to its optimality with respect to Hankelization and its somewhat standard use.

## 2.2. SSA-CP

The primary goal of this section is to present a simple method, SSA-CP, that improves the estimates of the final  $M-1$  entries of each RC, including in particular the current state estimate. In addition, SSA-CP will provide a prediction of RCs for  $t > N$ . (The same procedure may be directly applied to the first  $M-1$  entries of each reconstruction, but for simplicity of presentation, we focus solely on the last  $M-1$  entries.)

The steps of SSA-CP are as follows:

Step 1. Perform steps 1 and 2 of traditional SSA.

Step 2. Construct an extended lag-embedded matrix  $\vec{\mathbf{X}}$  of size  $(MD) \times N$ . The first  $N$  columns of  $\vec{\mathbf{X}}$  are identical to the columns of  $\mathbf{X}$ . For the final  $M-1$  columns, those entries that are known from the time series are filled in. The unknown entries below the diagonal consisting of  $x_N$ s are estimated using their conditional mean prediction,

$$\vec{\mathbf{X}} = \begin{bmatrix} \vec{x}_1 & \dots & \vec{x}_{N-M+1} & \vec{x}_{N-M+2} & \dots & \vec{x}_{N-1} & \vec{x}_N \\ \vec{x}_2 & \dots & \vec{x}_{N-M+2} & \vec{x}_{N-M+3} & \dots & \vec{x}_N & \vec{\mu}_{N+1|N} \\ \vdots & & \vdots & \vdots & & \vdots & \vdots \\ \vec{x}_{M-1} & \dots & \vec{x}_{N-1} & \vec{x}_N & \dots & \vec{\mu}_{N+M-3|N-1,N} & \vec{\mu}_{N+M-2|N} \\ \vec{x}_M & \dots & \vec{x}_N & \vec{\mu}_{N+1|N-M+2, \dots, N} & \dots & \vec{\mu}_{N+M-2|N-1,N} & \vec{\mu}_{N+M-1|N} \end{bmatrix}. \quad (6)$$

The calculation of each  $\vec{\mu}_{i|N-l, \dots, N}$  in (6) is as follows.

Let  $\vec{y}$  refer to the  $k$ th column of  $\vec{X}$ , with  $N + 1 \leq k \leq N + M - 1$ , and let  $\vec{y}_1, \vec{y}_2$  refer to the known and unknown portions of  $\vec{y} = [\vec{y}_1^T, \vec{y}_2^T]^T$ , respectively. If  $\vec{y}$  is a Gaussian random variable with mean  $\vec{\mu} = 0$  and covariance matrix  $\mathbf{C}$ , then  $\vec{y}_2$  has a conditional distribution that is Gaussian with mean

$$\vec{\mu}_{2|1} = \mathbf{C}_{21} \mathbf{C}_{11}^{-1} \vec{y}_1. \quad (7)$$

where  $\mathbf{C}$  can be written as

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{bmatrix} \quad (8)$$

with  $\mathbf{C}_{11}$  describing the covariance of the known values with themselves and so forth (Kaipio & Somersalo, 2005). The unknown entries  $\vec{y}_2$  are then filled in with the appropriate entries of  $\vec{\mu}_{2|1}$ , where  $\vec{\mu}_{N+j|k-M+1, \dots, N}$  in (6) denotes a  $D$ -dimensional column vector, that is, the  $j$ th set of  $D$  entries of the vector  $\vec{\mu}_{2|1}$ , calculated for the  $k$ th column of  $\vec{X}$  (with  $N + 1 \leq k \leq N + M - 1$ ). If necessary, a small amount of noise may be added to the covariance matrix in order to evaluate  $\mathbf{C}_{11}^{-1}$  in (7).

Step 3. Modify step 3 of traditional SSA by replacing  $\mathbf{X}$  with  $\vec{X}$ ; this change results in extended PCs  $\vec{\phi} = \vec{X}^T \vec{v}$ ; each extended PC is a column vector of length  $N$ .

Step 4. Modify step 4 of traditional SSA by replacing  $\phi$  with  $\vec{\phi}$  to construct an extended RC:

$$\vec{z}(t) = \frac{1}{\tilde{M}_t} \sum_{i=\tilde{L}_t}^{\tilde{U}_t} \vec{\phi}_{t-i+1} \vec{v}_i, \quad (9)$$

where  $(\tilde{M}_t, \tilde{L}_t, \tilde{U}_t)$  are defined by

$$(\tilde{M}_t, \tilde{L}_t, \tilde{U}_t) = \begin{cases} (t, 1, t), & 1 \leq t \leq M - 1, \\ (M, 1, M), & M \leq t \leq N, \\ (N - t + M, t - N + 1, M), & N + 1 \leq t \leq N + M - 1, \end{cases} \quad (10)$$

so that each extended RC  $\vec{z}$  is a (possibly multivariate) time series of length  $N + M - 1$ , with the last  $M - 1$  entries corresponding to predictions of the future state of the mode.

In the case that the data set has a Gaussian distribution, the conditional mean provides an optimal estimate of the missing data (Kaipio & Somersalo, 2005).

### 3. Data and Methods

The SSA-CP method will be tested on several data sets and compared to the traditional SSA reconstruction.

#### 3.1. Data

The first test uses a 15-year portion of the daily Real-time Multivariate MJO (RMM) indices (Wheeler & Hendon, 2004) from 1 January 1999 through 31 December 2013. The RMM indices have a distribution that is approximately normal with mean and variance approximately 0 and 1, respectively (Chen & Majda, 2015). For this two-dimensional data set,  $D = 2$  and  $N_{\text{tot}} = 5,479$ , with  $N_{\text{tot}}$  referring to the number of days.

Global Precipitation Climatology Project (GPCP) daily precipitation data (Huffman et al., 2012) are used for the second test. This data set has a spatial resolution of  $1^\circ \times 1^\circ$ ; the portion from 1 January 1997 through 31 December 2013 is used. Prior to applying SSA, the following steps were taken: (i) a meridional mode truncation to move from 2-D( $x, y$ ) to 1-D( $x$ ), (ii) removal of annual mean and seasonal cycle, and (iii) interpolation to 64 equally spaced zonal grid points. The meridional mode truncation step is a projection of the data onto the leading meridional mode proportional to  $e^{-y^2/2}$  where  $y$  is proportional to latitude; this step is identical to that used in, for example, Stechmann and Ogrosky (2014) and Stechmann and Majda (2015). Steps (i) and (iii) reduce the number of dimensions to  $D = 64$ , and the number of times is  $N_{\text{tot}} = 6,209$ . Note that these anomalies have a non-Gaussian distribution at each longitude; see the supporting information for the statistics of these anomalies.

A simulation of a multiscale model (Majda & Harlim, 2012) is used for the third test. The model equations are

$$du_1 = (-\gamma_1 u_1 + F(t)) dt + \sigma_1 dW_1, \quad (11a)$$

$$du_2 = (-\gamma_2 + i\omega_0/\epsilon + ia_0 u_1) u_2 dt + \sigma_2 dW_2, \quad (11b)$$

where  $\gamma_1 = \gamma_2 = 0.2$ ,  $\sigma_1 = \sigma_2 = 0.5$ ,  $\omega_0 = a_0 = 1$ ,  $\epsilon = 0.5$ , and  $F(t) = \sin(t/5)$ . An approximate solution was calculated numerically with the Euler-Maruyama method using  $dt = 0.005$  and  $t_{\text{end}} = 2,000$ . The real part of  $u_2$  was then sampled every 0.5 time units to create a data set with  $D = 1$  and  $N_{\text{tot}} = 4,000$ . A portion of this signal can be seen in Figure S2 in the supporting information.

### 3.2. Methods

The results of each real-time reconstruction method (SSA-CP and traditional) will be compared with the traditional reconstruction that has knowledge of future data. This is done in two steps.

First, both the traditional SSA and SSA-CP methods were applied to each data set after removing the final  $2M - 2$  time entries from the data set; for example, using an embedding window of  $M = 51$  days for the RMM indices, the methods were applied to the first  $N = N_{\text{tot}} - 2M + 2 = 5,379$  days. The embedding window was chosen to be large enough to be consistent with the intraseasonal time scale of the indices and is similar to that used in Chen and Majda (2015); other choices of this parameter value will be discussed in section 5. The standard reconstruction  $z(t)$  for each mode therefore has  $N = 5,379$  entries, while the SSA-CP reconstruction  $\tilde{z}(t)$  has  $N + M - 1 = 5,429$  entries. Note that the first  $N - M + 1 = 5,329$  entries for each reconstruction method are identical to one another; that is,  $z(t) = \tilde{z}(t)$  for  $1 \leq t \leq N - M + 1$ . Next, the traditional reconstruction method was used again, this time on the full  $N_{\text{tot}} = 5,479$  entries, resulting in a reconstruction  $u(t)$  with  $N_{\text{tot}} = 5,479$  entries. The entries of  $u(t)$  up to  $N_{\text{tot}} - M + 1 = 5,429$  are taken to be “truth,” and each of the methods applied to the shorter time series are compared with this truth.

Second, these tests are repeated for each data set with decreasing  $N_{\text{tot}}$ ; that is, define  $N_{\text{tot},i} = N_{\text{tot}} - i + 1$ , and repeat the test described above but using only the first  $N_{\text{tot}} = N_{\text{tot},i}$  entries of the data set, so that  $N = N_i := N_{\text{tot},i} - 2M + 2$ . For the RMM indices and multiscale model,  $i \in I = [1, \dots, 1001]$ ; for the GPCP data,  $i \in I = [1, 6, 11, \dots, 1001]$ . The pattern correlation and root-mean-square error (RMSE) are then calculated as a function of days before or after  $N_i$ ; see the supporting information for details.

## 4. Results

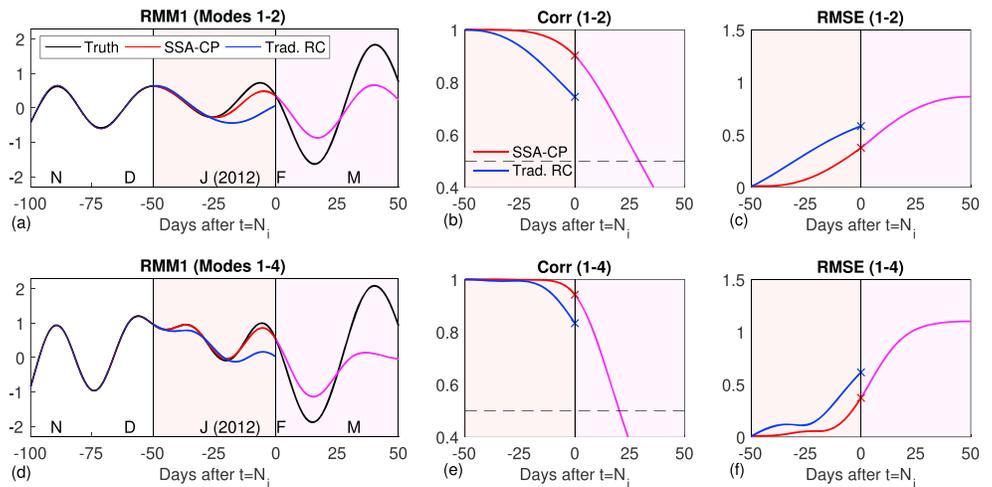
We next show results for three tests.

### 4.1. RMM Index

How well does the method perform on low-dimensional data that is nearly Gaussian?

Figures 1a and 1d show the results of using the SSA-CP or traditional reconstruction methods on the RMM indices with an embedding window  $M = 51$  days. For times away from the endpoints of the data, that is,  $t < N_i - M + 1$ , both methods are in agreement with the truth. For past times near the endpoints, that is,  $N_i - M + 1 < t < N_i$  (light orange-shaded region), SSA-CP captures both the phase and amplitude of the RMM1 index better than the traditional reconstruction. For future times  $t > N_i$ , SSA-CP is able to make predictions, with good agreement in phase and an underestimate of the amplitude of the true reconstruction. This underestimate of amplitude is due to using conditional mean predictions, which tend to 0 as  $t \rightarrow \infty$ .

The example in Figure 1 is a particularly challenging test as it is a case of Madden-Julian oscillation (MJO) onset. More specifically, the period being predicted in Figures 1a and 1d, namely, 1 February 2012 through 21 March 2012, exhibits a growing amplitude of the RMM1 index (black line), corresponding to the onset of the MJO event sometimes referred to as MJO4 during the 2011–2012 CINDY/DYNAMO field campaign (Yoneyama et al., 2013). This MJO has been considered a “primary” event, in that there are no clear signals connecting this MJO to the previous MJOs that occurred in October through December 2012. In contrast to the traditional reconstruction, SSA-CP more naturally captures oscillations and changes in frequency and amplitude, near the endpoint.

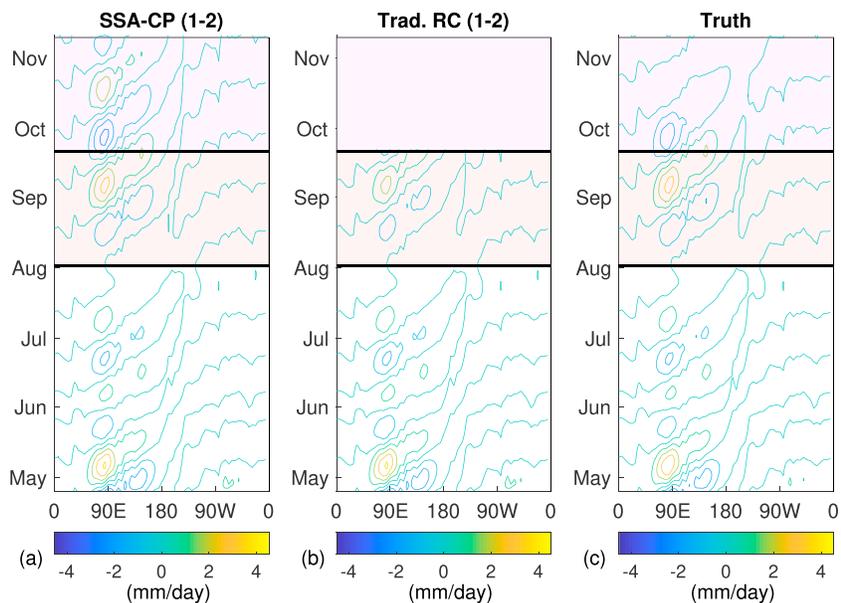


**Figure 1.** (a) Reconstructed RMM1 using components 1 and 2 with  $t = N_{601} = 4,779$  (31 January 2012) using (blue) traditional reconstruction, (red/magenta) SSA-CP, and (black) reconstruction using future information. (b, c) Bivariate pattern correlation and RMSE of the (blue) traditional reconstruction and truth as a function of days prior to/after  $N_i$ , and (red/magenta) SSA-CP reconstruction, and truth using modes 1 and 2. (d–f) Same as (a)–(c) but using components 1–4. RC = reconstructed component; RMM = Real-time Multivariate Madden-Julian oscillation; RMSE = root-mean-square error; SSA-CP = singular spectrum analysis with conditional predictions.

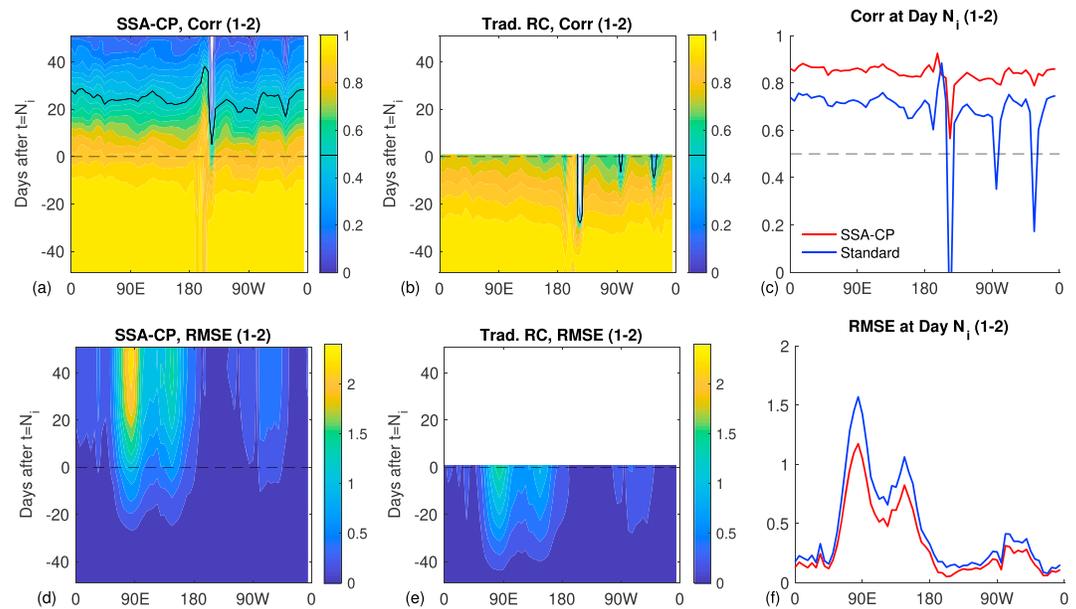
Figures 1b, 1c, 1e, and 1f show that when these tests are repeated, SSA-CP has significantly improved pattern correlation and reduced error compared to the traditional reconstruction. As a current state estimation, at  $t = N_i$ , SSA-CP improves the pattern correlation from 0.74 to 0.90 (0.83 to 0.94) for the two (four) leading modes. Likewise, SSA-CP reduces the error at  $t = N_i$  from 0.58 to 0.38 (0.62 to 0.37). For future times  $t > N_i$ , SSA-CP is able to make meaningful predictions for an extended period of time, with pattern correlations exceeding 0.5 out to approximately 29 (20) days when two (four) leading modes are used.

#### 4.2. Precipitation Data

How well does the method perform on large-dimensional, possibly non-Gaussian data?



**Figure 2.** (a) Reconstructed precipitation during 2013 using SSA-CP modes 1 and 2 with  $t_N = 6,109$ , corresponding to 22 September 2013. (b) Same as (a) but using traditional reconstruction. (c) Reconstructed modes 1 and 2 using future information. RC = reconstructed component; SSA-CP = singular spectrum analysis with conditional predictions.



**Figure 3.** (a) Pattern correlation, using 200 runs of SSA, of reconstructed precipitation components (1 and 2) using SSA-CP. (b) Same as (a) but for traditional reconstruction. (c) Pattern correlation at Day  $t = N_i$  for each method. (d–f) Same as (a)–(c) but showing RMSE. RC = reconstructed component; RMSE = root-mean-square error; SSA-CP = singular spectrum analysis with conditional predictions.

Figure 2 shows reconstructed precipitation anomalies using the two leading modes with an embedding window of 51 days. Both methods produce identical reconstructions prior to 2 August 2013. For 3 August 2013 through 22 September 2013, SSA-CP produces a reconstruction with amplitude in much better agreement with the non-real-time reconstruction (truth) than the traditional reconstruction. It also provides a prediction with decaying amplitude throughout October, qualitatively similar to the truth but with slower decay.

Repeating these tests for various  $N_i$  produces the pattern correlation and RMSE shown in Figure 3. For the recent past in time interval  $N_i - M + 1 < t < N_i$ , SSA-CP produces higher pattern correlation and lower RMSE than the standard reconstruction method. For state estimation at  $t = N_i$ , the pattern correlation is 0.1–0.2 higher at almost all longitudes when using SSA-CP than when using the standard method. Likewise, the RMSE is lower using SSA-CP than the traditional reconstruction at all longitudes. Note that low pattern correlation values for each method at longitudes like 150°W are due to small anomalies in the leading modes.

### 4.3. Partially Observed Multiscale Model

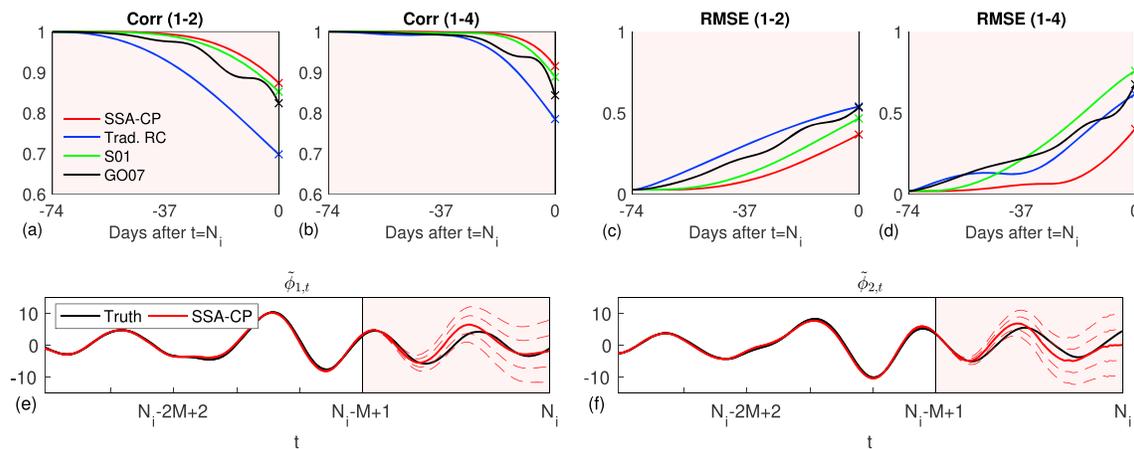
How well does the method perform on partially observed data?

Figure S3 in the supporting information shows the pattern correlation and RMSE for both methods applied to the multiscale model 11. For  $N_i - M + 1 < t < N_i$ , SSA-CP has significantly higher pattern correlation and lower error than the traditional reconstruction. At  $t = N_i$ , using SSA-CP improves the pattern correlation from 0.54 to 0.75 for two leading modes and lowers the error from 0.12 to 0.06. For  $t > N_i$ , predictions using SSA-CP have a pattern correlation of 0.5 or higher out to approximately 23 days when two leading modes are used.

## 5. Discussion

SSA-CP has been proposed as a method that supplements the mode identification ability of SSA with improved estimates of mode reconstructions near the ends of time series. We note that it is not at all necessarily the best possible data-driven, model-free prediction method that could be designed. Its effectiveness at identifying modes of variability in real time is of course also limited to cases where SSA is effective at identifying modes of interest.

How sensitive are the results to changes in the embedding window? As a first step toward addressing this question, the RMM tests from the previous section were rerun with an embedding window of  $M = 75$  days.



**Figure 4.** (a, c) Bivariate pattern correlation and RMSE of reconstructed Real-time Multivariate Madden-Julian oscillation indices using the (blue) traditional reconstruction, (red) SSA-CP reconstruction, (green) weighted reconstruction of Schoellhamer (2001), and (black)  $\Pi$  projector/simultaneous filling in method of Golyandina and Osipov (2007) as a function of days prior to  $N_i$ , using modes 1 and 2; here  $M = 75$ ,  $N_i = 4,779$ . (b, d) Same as (a) and (c) but for modes 1–4. (e, f) Leading two principal components of SSA-CP. Dashed red lines indicate  $\pm 1, 2$  standard deviations. RC = reconstructed component; RMSE = root-mean-square error; SSA-CP = singular spectrum analysis with conditional predictions.

Figures 4a–4d show that while both SSA-CP and the traditional reconstruction produce slightly lower pattern correlation at  $t = N_i$  than in the previous test with  $M = 51$ , SSA-CP again results in significantly higher PC and lower RMSE than the traditional reconstruction. For  $t > N_i$ , the pattern correlation stays higher than 0.5 for 35 (24) days when the leading two (four) modes are used (not shown). Results of additional tests using  $M = 61, 71, 81$ , and 101 are shown in Figure S4 in the supporting information; increasing the embedding window provides some small improvement in the pattern correlation and RMSE of predictions, which is likely due to the increased time scales present in the modes identified with a larger embedding window. Further tests indicate that, in addition, the method performs equally well when the leading modes are found using a period of training data that is not near the endpoint; see Figure S5 in the supporting information for further details.

How does SSA-CP compare with other methods in the literature that have been proposed for either (i) improving state estimation of RCs near the endpoints of time series or (ii) using SSA on data sets with gaps? We briefly examine this through a comparison of the results of SSA-CP with methods from Schoellhamer (2001) and Golyandina and Osipov (2007) for the first test from section 4. Figures 4a–4d show the pattern correlation and RMSE of these two methods along with the traditional reconstruction and SSA-CP. All of the modified versions of SSA produce higher pattern correlation than the traditional reconstruction, with SSA-CP having the highest. For the leading two modes, all methods produce lower RMSE than the traditional reconstruction, but when the leading four modes are used, only SSA-CP outperforms the traditional reconstruction over each of the final  $M - 1$  days. Other methods for using SSA (or other mode identification methods) in real-time have been proposed, and it would be interesting to investigate further comparisons in a future study. For example, other methods include a predicted spatial basis method (Chen et al., 2018), kernel analog forecasting (e.g., Comeau et al., 2018), methods based on linear recurrent formulas (Golyandina et al., 2001), methods that project smoothed data onto leading SSA modes computed with Fourier-filtered data (Roundy & Schreck III, 2009), and energy-minimizing reconstructions of PCs (Shen et al., 2015).

Many additional tests were conducted beyond the three examples described in detail above. Other tests were conducted using data sets generated by stochastic processes (complex-valued Ornstein-Uhlenbeck process), deterministic dynamical systems (Lorenz 63 model, multiple examples from Golyandina et al., 2001), other observational data (Kelvin wave calculated using National Centers for Environmental Prediction/National Center for Atmospheric Research reanalysis data [Kalnay et al., 1996] and the methods of Ogrosky and Stechmann (2015, 2016), and numerous synthetic test signals both with and without noise. SSA-CP significantly outperformed traditional SSA in almost all of these tests. In cases of deterministic signals of Golyandina et al. (2001), both methods produced excellent reconstructions of the leading modes near the endpoints. In cases like this, the standard reconstruction may be just as desirable as SSA-CP or any other modification, as the additional effort of implementing SSA-CP, though minimal, may not be necessary to provide reasonable

initial conditions for a forecast. In addition, one benefit of the standard reconstruction is its invertibility; if all modes are reconstructed and summed together, the original data set is recovered. This invertibility is not shared by SSA-CP.

There are several compelling reasons for using SSA-CP rather than the traditional reconstruction, however. First, it is nearly as simple to use as traditional SSA. Second, it is optimal for Gaussian data and is based on well-known theory. Third, it is straightforward to quantify the uncertainty in the extended PCs or reconstruction. For example, the variance of  $\tilde{\phi}_{N-l+1}$ , where  $1 \leq l \leq M - 1$ , is given by

$$\text{Var}(\tilde{\phi}_{N-l+1}) = [\tilde{v}_{l+1}^T, \dots, \tilde{v}_M^T] \mathbf{C}_{2l} [\tilde{v}_{l+1}^T, \dots, \tilde{v}_M^T]^T. \quad (12)$$

Figures 4e and 4f show the two leading PCs of the RMM indices calculated using SSA-CP with  $N_l = 4, 779$  and  $M = 75$ . One and two standard deviations from the extended PC entries are shown, with the standard deviation calculated using (12).

Finally, since non-Gaussianity leads to a lack of independence between modes in linear methods like EOFs, there is no guarantee that the method will work well on data with strong non-Gaussianity (Monehan et al., 2009). However, the method works well on the non-Gaussian data used here, perhaps owing to the somewhat mild deviations from Gaussianity. The method could potentially be extended to non-Gaussian frameworks with conditional Gaussian or Gaussian mixture structures (see, e.g., Chen & Majda, 2018; Majda, 2016).

We note that SSA is just one of many data analysis tools capable of identifying modes of variability in spatiotemporal data sets (see Crommelin & Majda, 2004, for a discussion of other linear methods for mode identification). SSA was chosen to be the focus of the current study due to its linearity, simplicity, and popularity, combined with the linearity of the proposed modifications. Other mode identification methods, including nonlinear methods like nonlinear Laplacian spectral analysis, have been shown to be effective at capturing modes of variability that SSA has difficulty capturing, like modes with pronounced intermittent behavior (Giannakis & Majda, 2012a, 2012b) and theory supporting both such methods and forecasting techniques of relevance has been developed in recent years (Comeau et al., 2018; Zhao & Giannakis, 2016). Including conditional predictions into such methods is certainly possible; methods like nonlinear Laplacian spectral analysis use a reconstruction approach similar to that of (4) and incorporating conditional predictions into this method could potentially be done in a straightforward manner. It is not clear, however, that such a method would be optimal in the same way that conditional predictions used here are optimal when combined with Gaussian data and with a linear method like SSA.

## 6. Conclusions

In summary, a modified SSA algorithm, SSA-CP, has been presented and tested. This modification is proposed to address endpoint issues that arise when using SSA. When compared with the traditional reconstruction method, SSA-CP results in significantly improved real-time estimates of leading modes of variability when applied to a variety of data sets.

This method was shown to be useful for providing improved initial conditions for forecasts. It is derived from well-known theory using Gaussian statistics and provides optimal predictions for Gaussian data, but also performs well in tests with non-Gaussian data. The uncertainty in the real-time estimates may be quantified using the covariance matrix that is inherently part of the method.

While the current study has been primarily focused on applying the method to atmospheric science data, this method may prove useful in application areas outside of atmospheric science. In addition, it is possible that the ideas used here may be adapted for other methods of mode identification. These subjects are left for future work.

## References

- Aubry, N., Guyonnet, R., & Lima, R. (1991). Spatiotemporal analysis of complex signals: Theory and applications. *Journal of Statistical Physics*, *64*, 683–739.
- Broomhead, D. S., & King, G. P. (1986). Extracting qualitative dynamics from experimental data. *Physica D*, *20*, 217–236.
- Chen, N., & Majda, A. J. (2015). Predicting the real-time multivariate Madden-Julian oscillation index through a low-order nonlinear stochastic model. *Monthly Weather Review*, *143*, 2148–2169.

### Acknowledgments

The GPCP data for this article are available from NOAA/OAR/ESRL PSD, Boulder, Colorado, USA, from their website (<http://www.esrl.noaa.gov/psd/>). The RMM indices can be obtained online at Bureau of Meteorology website (<http://www.bom.gov.au/climate/mjo/>). Other data used are in the figures. The research of S. N. S. is partially supported by a Sloan Research Fellowship from the Alfred P. Sloan Foundation and a Vilas Associates Award from the University of Wisconsin-Madison. The research of N. C. is supported by the Office of Vice Chancellor for Research and Graduate Education (VCRGE) at University of Wisconsin-Madison.

- Chen, N., & Majda, A. J. (2016). Filtering the stochastic skeleton model for the Madden-Julian oscillation. *Monthly Weather Review*, *144*, 501–527.
- Chen, N., & Majda, A. J. (2018). Conditional Gaussian systems for multiscale nonlinear stochastic systems: Prediction, state estimation and uncertainty quantification. *Entropy*, *20*, 509.
- Chen, N., Majda, A. J., Sabeerali, C. T., & Ravindran, A. J. (2018). Predicting monsoon intraseasonal precipitation using a low-order stochastic model. *Journal of Climate*, *31*, 4403–4427.
- Comeau, D., Giannakis, D., Zhao, Z., & Majda, A. J. (2018). Predicting regional and pan-Arctic sea ice anomalies with kernel analog forecasting. *Climate Dynamics*. <https://doi.org/10.1007/s00382-018-4459-x>
- Crommelin, D. T., & Majda, A. J. (2004). Strategies for model reduction: Comparing different optimal bases. *Journal of the Atmospheric Sciences*, *61*, 2206–2217.
- Ghil, M., Allen, R. M., Dettinger, M. D., Ide, K., Kondrashov, D., Mann, M. E., et al. (2002). Advanced spectral methods for climatic time series. *Reviews of Geophysics*, *40*(1), 1003. <https://doi.org/10.1029/2000RG000092>
- Giannakis, D., & Majda, A. J. (2012a). Nonlinear Laplacian spectral analysis for time series with intermittency and low-frequency variability. *Proceedings of the National Academy of Sciences of the United States of America*, *109*, 2222–2227.
- Giannakis, D., & Majda, A. J. (2012b). Comparing low-frequency and intermittent variability in comprehensive climate models through nonlinear Laplacian spectral analysis. *Geophysical Research Letters*, *39*, L10710. <https://doi.org/10.1029/2012GL051575>
- Golyandina, N., Nekrutkin, V., & Zhigljavsky, A. A. (2001). Analysis of time series structure. SSA and Related Techniques, Chapman and Hall, CRC.
- Golyandina, N., & Osipov, E. (2007). The “Caterpillar”-SSA method for analysis of time series with missing values. *Journal of Statistical Planning and Inference*, *137*, 2642–2653.
- Hassani, H. (2007). Singular spectrum analysis: Methodology and comparison. *Journal of Data Science*, *5*, 239–257.
- Hassani, H., Webster, A., Silva, E. S., & Heravi, S. (2014). Forecasting U.S. tourist arrivals using optimal singular spectrum analysis. *Tourism Management*, *46*, 322–335.
- Huffman, G. J., Bolvin, D. T., & Adler, R. F. (2012). GPCP version 2.2 SG combined precipitation data set. WDC-A, NCDC, Asheville, NC. Data set accessed 12 February 2014 at <http://www.ncdc.noaa.gov/oa/wmo/wdcamet-ncdc.html>
- Kaipio, J., & Somersalo, E. (2005). *Statistical and computational inverse problems*. New York: Springer.
- Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., et al. (1996). The NCEP/NCAR 40-year reanalysis project. *Bulletin of the American Meteorological Society*, *77*, 437–471.
- Kang, I.-S., & Kim, H.-M. (2010). Assessment of MJO predictability for boreal winter with various statistical and dynamical models. *Journal of Climate*, *23*, 2368–2378.
- Keppenne, C. L., & Ghil, M. (1990). Adaptive filtering and prediction of the Southern Oscillation Index. *Journal of Geophysical Research*, *97*, 20,449–20,454.
- Kikuchi, K., & Wang, B. (2008). Diurnal precipitation regimes in the global tropics. *Journal of Climate*, *21*, 2680–2696.
- Kondrashov, D., Chekroun, M. D., Robertson, A. W., & Ghil, M. (2013). Low-order stochastic model and “past-noise forecasting” of the Madden-Julian oscillation. *Geophysical Research Letters*, *40*, 5305–5310. <https://doi.org/10.1002/grl.50991>
- Kondrashov, D., & Ghil, M. (2006). Spatio-temporal filling of missing points in geophysical data sets. *Nonlinear Processes in Geophysics*, *13*, 151–159.
- Kondrashov, D., Shprits, Y., & Ghil, M. (2010). Gap filling of solar wind data by singular spectrum analysis. *Geophysical Research Letters*, *37*, L15101. <https://doi.org/10.1029/2010GL044138>
- Lisi, F., & Medio, A. (1997). Is a random walk the best exchange rate predictor? *International Journal of Forecasting*, *13*, 255–267.
- Majda, A. J. (2016). *Introduction to turbulent dynamical systems in complex systems*. Cham, Switzerland: Springer.
- Majda, A. J., & Harlim, J. (2012). *Filtering complex turbulent systems*. New York: Cambridge University Press.
- Mo, K. C. (2001). Adaptive filtering and prediction of intraseasonal oscillations. *Monthly Weather Review*, *129*, 802–817.
- Monehan, A. H., Frye, J. C., Ambaum, M. H. P., Stephenson, D. B., & North, G. R. (2009). Empirical orthogonal functions: The medium is the message. *Journal of Climate*, *22*, 6501–6514.
- Ogrosky, H. R., & Stechmann, S. N. (2015). Assessing the equatorial long-wave approximation: Asymptotics and observational data analysis. *Journal of the Atmospheric Sciences*, *72*, 4821–4843.
- Ogrosky, H. R., & Stechmann, S. N. (2016). Identifying convectively coupled equatorial waves using theoretical wave eigenvectors. *Monthly Weather Review*, *144*, 2235–2264.
- Rodrigues, P. C., & de Carvalho, M. (2013). Spectral modeling of time series with missing data. *Applied Mathematical Modelling*, *37*, 4676–4684.
- Roundy, P. E., & Schreck III, C. J. (2009). A combined wave-number–frequency and time-extended EOF approach for tracking the progress of modes of large-scale organized tropical convection. *Quarterly Journal of the Royal Meteorological Society*, *135*, 161–173.
- Schoellhamer, D. H. (2001). Singular spectrum analysis for time series with missing data. *Geophysical Research Letters*, *28*, 3187–3190.
- Shen, Y., Li, W., Xu, G., & Li, B. (2014). Spatiotemporal filtering of regional GNSS. *Journal of Geodesy*, *88*, 1–12.
- Shen, Y., Peng, F., & Li, B. (2015). Improved singular spectrum analysis for time series with missing data. *Nonlinear Processes in Geophysics*, *22*, 371–376.
- Stechmann, S. N., & Majda, A. J. (2015). Identifying the skeleton of the Madden-Julian oscillation in observational data. *Monthly Weather Review*, *143*, 395–416.
- Stechmann, S. N., & Ogrosky, H. R. (2014). The Walker circulation, diabatic heating, and outgoing longwave radiation. *Geophysical Research Letters*, *41*, 9097–9105. <https://doi.org/10.1002/2014GL062257>
- Vautard, R., & Ghil, M. (1989). Singular spectrum analysis in non-linear dynamics, with applications to paleoclimatic time series. *Physica D*, *35*, 395–424.
- Vautard, R., Yiou, P., & Ghil, M. (1992). Singular spectrum analysis: A toolkit for short noisy chaotic signals. *Physica D*, *58*, 95–126.
- Weare, B. C., & Nasstrom, J. S. (1982). Examples of extended empirical orthogonal function analyses. *Monthly Weather Review*, *110*, 481–485.
- Wheeler, M. C., & Hendon, H. H. (2004). An all-season real-time multivariate MJO index: Development of an index for monitoring and prediction. *Monthly Weather Review*, *132*, 1917–1932.
- Yoneyama, K., Zhang, C., & Long, C. N. (2013). Tracking pulses of the Madden-Julian oscillation. *Bulletin of the American Meteorological Society*, *94*, 1871–1891.
- Zhao, Z., & Giannakis, D. (2016). Analog forecasting with dynamics-adapted kernels. *Nonlinearity*, *29*, 2888–2939.