

An empirical stochastic model of sea-surface temperatures and surface winds over the Southern Ocean

S. Kravtsov¹, D. Kondrashov², I. Kamenkovich³, and M. Ghil^{2,4}

¹University of Wisconsin-Milwaukee, Dept. of Mathematical Sciences, Atmospheric Science group, P.O. Box 413, Milwaukee, WI 53201, USA

²University of California at Los Angeles, USA

³University of Miami, USA

⁴Ecole Normale Supérieure, Paris, France

Received: 1 August 2011 – Published in Ocean Sci. Discuss.: 29 August 2011

Revised: 26 October 2011 – Accepted: 2 November 2011 – Published: 14 November 2011

Abstract. This study employs NASA's recent satellite measurements of sea-surface temperatures (SSTs) and sea-level winds (SLWs) with missing data filled-in by Singular Spectrum Analysis (SSA), to construct empirical models that capture both intrinsic and SST-dependent aspects of SLW variability. The model construction methodology uses a number of algorithmic innovations that are essential in providing stable estimates of the model's propagator. The best model tested herein is able to faithfully represent the time scales and spatial patterns of anomalies associated with a number of distinct processes. These processes range from the daily synoptic variability to interannual signals presumably associated with oceanic or coupled dynamics. Comparing the simulations of an SLW model forced by the observed SST anomalies with the simulations of an SLW-only model provides preliminary evidence for the ocean driving the atmosphere in the Southern Ocean region.

SST-dependent SLW variability can help analyze the coupled climate dynamics of the Southern Ocean, especially when combined with oceanic General Circulation Models (GCMs).

Climate variability over the Southern Ocean is likely to be of global significance, due to this ocean's special role in linking the Atlantic, Pacific, and Indian basins. However, progress in understanding the dynamics of large-scale air-sea coupling over the Southern Ocean has been slow, largely due to the very low density of in situ measurements in this region. Recently launched NASA satellites provide accurate high-resolution global measurements of important climatic variables such as SST and SLW. These global fields now permit the construction of empirical air-sea interaction models for the Southern Ocean. Despite improved data coverage in the region, estimating the propagator of the above mentioned statistical models remains an ambitious and challenging task, since (1) there are still missing data due to the presence of strong winds or heavy rains, and (2) such a model has to have an unprecedentedly large number of degrees of freedom, due to high-dimensional nature of global-scale air-sea interaction. The model construction thus requires major algorithmic revisions and gap-free datasets, which we develop and describe in detail below.

1 Introduction

1.1 Motivation

This study addresses aspects of ocean-atmosphere interaction over the Southern Ocean using measurements provided by satellite sensors. Our objective is to quantitatively describe and analyze co-variability of sea-surface temperature (SST) and sea-level wind (SLW) in this region, by developing inverse stochastic models that are derived directly from the remotely sensed data. Empirical models of potentially

1.2 Background

The Southern Ocean is the region south of roughly 30° S that includes the Antarctic Circumpolar Current (ACC), along with the branches of circulations that link it to the Atlantic, Pacific, and Indian Oceans (Schmitz, 1996).



Correspondence to: S. Kravtsov
(kravtsov@uwm.edu)

1.2.1 Satellite data over the Southern Ocean

Poor spatial coverage by in situ measurements in the Southern Ocean prohibits direct comprehensive description of climate variability there. The data from NASA satellites launched over the past decade thus provide a unique source of precise measurements of climatically important quantities such as SST and SLW. Global coverage and fine resolution make them extremely valuable for studying air-sea interaction in the Southern Ocean. In particular, the microwave-based sensor AMSR on the AQUA satellite launched in 2002 samples SST field under clouds – an opportunity that was previously unavailable for infrared-based SST records over the typically cloudy Southern Ocean.

Microwave-based SST products (Kummerow et al., 2000; Wentz et al., 2000) have been utilized before to explore tropical SST variations (Hashizume et al., 2000; Chelton et al., 2001; Harrison and Vecchi, 2001; Vecchi and Harrison, 2002; Vecchi et al., 2003). We will use this type of measurements in the present paper to address air-sea interaction over the Southern Ocean.

1.2.2 Southern Ocean climate

The Southern Ocean is characterized by intense climatological westerlies that induce strong meridional Ekman transports and drive the ACC. The modes of climate variability in the Southern Ocean differ by their time scales and spatial signatures, as well as by specific dynamical mechanisms. Synoptic variability, with a time scale of a few days, is comprised of extremely powerful atmospheric storms associated with baroclinic Rossby waves passing over the region. These synoptic eddies cause large-amplitude SST responses mainly through enhanced vertical turbulent mixing in the oceanic boundary layer and through changes in the air-sea heat flux. An example of such coherent patterns of wind and SST anomalies associated with synoptic variability is shown in Fig. 1, which displays snapshots of these fields' anomalies on 1 December 2002; the anomalies were computed relative to the base 16-day period of 1–16 December 2002. The pattern correlation between the SST and SLW fields over the region shown in Fig. 1 is of about $r = -0.73$, which allows to reject the null hypothesis of zero correlation in favor of the alternative of negative correlation at 0.1 % level, according to the one-sided t -test with $\nu = 16$ degrees of freedom (the value of t statistic is $t = r\sqrt{\nu/(1-r^2)} \approx -4.3$). The number of spatial degrees of freedom within the region of interest in the above test was estimated based on the decorrelation scale of about 1000 km.

On longer time scales, an intraseasonal mode of intrinsic atmospheric variability is called the Southern Annular Mode (SAM). It has a pronounced zonally symmetric component, hence its name (Thompson and Wallace, 2000; Thompson et al., 2000). It is also known as zonal-flow vacillation (Hartmann, 1995) and consists of irregular meridional dis-

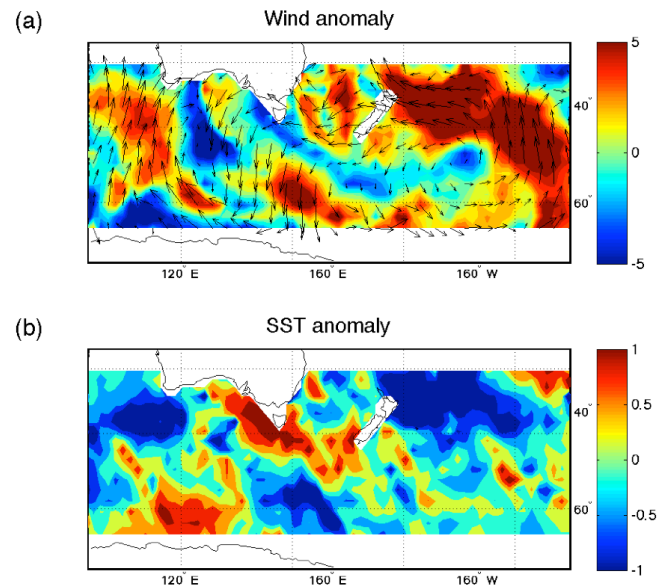


Fig. 1. Anomalies of atmospheric (SLW) and oceanic (SST) fields on 1 December 2002, computed as the deviations from the average over the 16-day period of 1–16 December 2002: (a) sea-level wind (SLW) anomaly; (b) sea-surface temperature (SST) anomalies. The two fields are spatially correlated, over the region shown, with the correlation coefficient of -0.73 .

placements of the atmospheric jet. SAM is thought to be energized by higher-frequency synoptic eddies and may, in turn, modify the storm track at lower frequencies (Robinson, 2000; Lorenz and Hartmann, 2001). Feldstein (2000) has argued that SAM variability is due to linear dynamical response to stochastic forcing associated with synoptic eddies. In contrast, Koo et al. (2002) presented a non-linear framework for zonal-flow vacillation, based on the paradigm of weather regimes (Reinhold and Pierrehumbert, 1982; Legras and Ghil, 1985; Marshall and Molteni, 1993; Koo and Ghil, 2002; Kravtsov et al., 2005a). Hall and Visbeck (2002) discussed the dynamics of oceanic response to SAM-type surface-wind evolution, and reported significant variations in SST, sea-ice extent and ACC transport associated with this variability. Even longer-term modes of variability include the so-called semi-annual oscillation (Van Loon, 1967, 1972; Meehl, 1991; Meehl et al., 1998) and the Pacific-South American (PSA) oscillation (Mo and Ghil, 1987) often described as a tropically forced standing wave train (Mo and White, 1985; Mo and Ghil, 1987; Karoly, 1989; Grimm and Silva Dias, 1995; Garreaud and Battisti, 1999). The PSA has also been associated with an El Niño/Southern Oscillation (ENSO) teleconnection pattern (Kwok and Comiso, 2002) and has signatures in the Southern Ocean's surface air temperature, SST and sea-ice extent; these signatures are referred to as the Antarctic Dipole (Yuan and Martinson, 2000, 2001). Finally, the Antarctic Circumpolar Wave (ACW) has an interannual time scale and

is associated with eastward-propagating signals in SST, sea-level pressure, and sea-ice extent (White and Peterson, 1996; Jacobs and Mitchell, 1996; Peterson and White, 1998).

1.2.3 Air-sea coupling over the Southern Ocean

Synoptic eddies and the lower-frequency SAM that dominate climate variability in the Southern Ocean on weekly-to-intraseasonal time scales are due primarily to intrinsic atmospheric dynamics. Surface manifestations of these modes induce significant oceanic response. The SST variability associated with this response affects, in turn, the atmospheric flow. Additionally, the SLW variability may be modified, on an intraseasonal-to-interannual and longer time scale, by SST anomalies associated with intrinsic oceanic or inherently coupled processes, such as the PSA and ACW.

Surface wind influences SST directly by modifying vertical turbulent heat exchange between the two fluids (Gill, 1982; Arya, 1988) and inducing strong horizontal Ekman transports in the oceanic mixed layer. High-frequency wind forcing also leads to significant long-term oceanic changes by affecting, among other things, the seasonal-mean subsurface temperatures and mixed-layer depths (Kamenkovich, 2005). In addition, surface-wind fluctuations can energize intrinsic oceanic modes, which may play an important role in the dynamics of the Southern Ocean (Wunsch, 1999; Weisse et al., 1999; Karsten et al., 2002; Gille, 2003). The SST signatures of these modes have structures that are different from that of a local SST response to wind forcing. These oceanic phenomena have long intrinsic time scales and may thus lead to partial predictability of the Southern Ocean climate.

The way SLW may respond to SST anomalies is via changes in stability of the marine atmospheric boundary layer. Air passing over a positive SST anomaly becomes more unstable; this leads to anomalous turbulent momentum flux and amplification of the surface wind (Arya, 1988). This effect was shown to be at work over the Eastern Tropical Pacific (Wallace et al., 1989; Liu et al., 2000; Chelton et al., 2001; Hashizume et al., 2001) and over the Southern Ocean (O'Neill et al., 2003) on seasonal-to-interannual time scales. Other dynamical factors may also contribute to this response at all time scales (Hsu, 1984; Lindzen and Nigham, 1987; Mitchell and Wallace, 1992).

The SST-induced modifications of the atmospheric boundary layer may cause changes in the free atmosphere's circulation. The linear response is expected to be weak, but nonlinear modes of atmospheric variability, such as SAM, may produce a stronger effect (Koo et al., 2002; Kravtsov et al., 2006a, b). Feliks et al. (2004, 2007) have shown, in particular, how an oceanic thermal front may induce intraseasonal variability in the overlying atmosphere, including surface-wind evolution.

To summarize, the Southern Ocean is characterized by vigorous variability on a wide range of time scales. Air-sea interaction in the region is complex and difficult to repre-

sent in dynamical models, as it involves a wide variety of boundary layer processes, as well as their coupling to intrinsic dynamics of the fluids on both sides of the ocean-atmosphere interface. Statistically, however, these interactions may well be described by joint variability of SLW and SST. A purely empirical model of this co-variability, based on recent high-quality satellite observations, could provide an accurate quantitative description of air-sea interaction without having to resolve explicitly the complex chain of participating dynamical processes.

1.2.4 Empirical stochastic models of SST and SLW

Data-based inverse stochastic models used in climate dynamics generally belong to one of the two major groups: (i) multivariate parametric models with additive, state-independent noise, the simplest of which is the so-called linear inverse model (LIM) (Penland, 1989, 1996; Penland and Sardeshmukh, 1995; Penland and Matrosova, 1998; Winkler et al., 2001); and (ii) nonparametric, univariate or bivariate models involving state-dependent, multiplicative noise (Sura, 2003; Sura and Gille, 2003; Sura et al., 2006; Sura and Newman, 2008; Sura and Sardeshmukh, 2008). Both types of models can be useful in addressing various aspects of climate variability, but are very different in terms of how they are constructed, as well as in their potential applications.

In particular, the models with multiplicative noise consider the time series of a variable of interest (for example, u -component of surface wind, or SST) at a single spatial location, and estimate state-dependent drift and diffusion parameters of the stochastic differential equation (SDE), which presumably governs the evolution of this variable. In order to get reliable estimates of model parameters given relatively sparse observations, as is the case for Southern Ocean winds, one may concatenate data sets from multiple locations, which are situated far enough so that their respective time series may be assumed to be uncorrelated (Sura, 2003). The scalar SDEs so obtained describe local features of interactions between processes evolving on different time scales. They are particularly successful in interpreting some of the nongaussian aspects of both SLW (Sura, 2003; Monahan, 2004, 2006a,b) and SST (Sura et al., 2006; Sura and Newman, 2008; Sura and Sardeshmukh, 2008) variability. The multiplicative noise in the above studies is attributed to random fluctuations of the drag coefficient or air-sea heat exchange coefficient.

On the other hand, multivariate parametric models driven by additive noise are usually constructed in the phase space of the leading empirical orthogonal functions (EOFs) (Preisendorfer, 1988) of the field(s) of interest, thus addressing non-local aspects of the variability under consideration. This non-locality comes at the expense of a fairly restrictive parametric dependency of the system's tendency on its state. In LIMs, for example, this dependency is assumed to be linear, while the model coefficients and noise parameters are

found by multiple linear regression (MLR). LIMs driven by Gaussian stochastic forcing cannot model the nongaussian aspects of the observed statistics, but more general, nonlinear empirical parametric models can. Kravtsov et al. (2005b) developed a methodology for constructing such nonlinear empirical models, which also addresses some other weaknesses of LIMs. This methodology showed excellent results when applied to the problems of mid-latitude variability of geopotential heights (Kondrashov et al., 2006), as well as to describing tropical SST evolution (Kondrashov et al., 2005).

1.3 This paper

The purpose of the present paper is to construct an empirical model of SLW variability over the Southern Ocean by using concurrent high-quality satellite measurements of SLW and SST. Doing so requires the use of recent microwave-sensed SST fields available after the launch of AQUA in June 2002. Since only about 5 years of such data are available, we do not attempt to develop a closed model that would simulate by itself long-term aspects of SLW-SST co-variability, such as ACW; this would require a much longer data set with enough degrees of freedom to capture interannual SST signals. Instead, the quantities involving SST observations will serve as predictors in the stochastic model of SLW evolution; the time-dependent SST anomalies themselves will be treated as given. We will show that this model is capable of reproducing the statistics of daily-to-intraseasonal SLW anomalies. As a brief introductory illustration of one of many potential uses of the empirical model constructed, we will present some evidence for large-scale oceanic imprint onto the atmospheric variability in the Southern Ocean by comparing the statistics of an SLW-only empirical model with that of a model forced by the daily history of SST anomalies.

Our statistical SST-dependent SLW model will also be able to capture some aspects of air-sea interaction and longer-term variability when coupled to a dynamical oceanic component. Experiments with such a coupled dynamical-statistical model will be studied in a future paper. The application of our statistical model as a component of a hybrid coupled GCM requires that both local and non-local aspects of SLW variability and its coupling with SST variability be comprehensively represented in the empirical model. We will therefore build upon the methodology of Kravtsov et al. (2005b) to construct this model, but emphasize here that substantial modifications to that model construction technique are necessary, as detailed below.

As we have mentioned at the end of section 1.1, the high-dimensional nature of basin-scale air-sea coupling in the Southern Ocean region prohibits direct application of Kravtsov et al. (2005b) method and requires major modifications to the model construction algorithm; these changes, when applied to gap-free satellite datasets with missing data filled-in by M-SSA (Kondrashov and Ghil, 2006), are essen-

tial in obtaining stable estimates of the empirical model propagator.

The rest of the paper is organized as follows. Section 2 describes the data sources, pre-processing and gap-filling methodology, as well as the data set's basic statistics. Section 3 outlines general, as well as novel technical aspects of the empirical stochastic model construction, with methodological details given in the appendices. The performance of our empirical models is evaluated in Sect. 4, while Sect. 5 summarizes our results and elaborates on their significance.

2 Data, pre-processing methodology, and basic statistics

2.1 Data sources

The gridded data products used in this analysis are obtained from the Remote Sensing Systems website (<http://www.ssmi.com>). The SST data are taken from the AMSR-E ocean data product (Version-5) for the time interval from June 2002 to February 2007 (Kawanishi et al., 2003). Missing data are due to sun glint, heavy rain, proximity of ice edge, and winds greater than 20 m s^{-1} . The wind speed and direction at 10 m a.s.l. are obtained from the QuikSCAT scatterometer dataset (Liu, 2002). The geophysical data record began on July 1999; for the analysis in this paper, we use data for the time interval that overlaps with that of the AMSR-E dataset. Although the scatterometer data tend to be less accurate in the presence of rain, we do not remove such data entries, since our statistical technique is based on analyzing spatial covariances within the fields considered; the small-scale random errors associated with rain occurrences will thus be effectively filtered out.

Both gridded data sets are available on a $0.25^\circ \times 0.25^\circ$ grid twice a day, on ascending and descending paths. For our subsequent analyses, the data were averaged in space and time to produce daily values on a $2^\circ \times 2^\circ$ grid.

2.2 Filling the missing data

While recent satellite observations over the Southern Ocean do have a previously unprecedented quality, there are still gaps in data coverage in the presence of heavy rains or strong winds; specifically, about 40 % of the points in the SST data set and 20 % in the SLW data set were missing. In order to fill these gaps in the data, we used the methodology of Kondrashov and Ghil (2006). Their algorithm is based on multi-channel singular spectrum analysis (M-SSA) (Ghil et al., 2002) and takes advantage of both spatial and temporal correlations in the existing data to iteratively produce estimates of missing data points, which are then used to compute a self-consistent spatiotemporal lag-covariance matrix; cross-validation is applied to find the optimal window width and number of dominant M-SSA modes to fill the gaps.

The missing data have been filled-in for SLW and SST fields separately; that is, cross-correlations between these

two fields were not exploited. Since the total number of spatial grid points exceeds the temporal length of the data set (in days) for both SLW and SST, we utilized the “reduced-covariance” approach (Ghil et al., 2002) to compute the spatio-temporal lag-covariance matrix. Based on the results of cross-validation experiments, we chose the lag of 1 day and 300 M-SSA modes for filling SLW components and the lag of 5 days and 160 M-SSA modes for SST. The domain-averaged root-mean-square (rms) error for filled-in values is estimated to be 0.44°C for SSTs and 1.7 m s^{-1} for SLW components.

2.3 Basic statistics

2.3.1 Filtering

We considered continuous, filled-in by M-SSA data sets of daily SST scalars and SLW vectors on a $2^{\circ} \times 2^{\circ}$ grid ($65^{\circ} - 30^{\circ}\text{S}$), for the period of 1 June 2002–13 February 2007, for a total of 1719 days. We first removed the seasonal cycle by retaining, at each grid point, only the residual of the multiple linear regression of the original, unfiltered time series onto a ten-variable “seasonal cycle” time series. The latter time series had the form $(\sin(2\pi nt/365), \cos(2\pi nt/365))$, where time t is measured in days and changes from $t = 0$ to $t = 1718$, while $n = 1, 2, 3, 4, 5$. The filtered versions of the original SST and SLW time series were then also linearly detrended to get rid of secular variability, since our statistical models are assumed to be stationary.

2.3.2 Low-order moments

Figure 2 shows a few of low-order moments of the filtered anomalies so obtained. The time-mean wind is plotted in panel (a), with the wind speed given by color shading, and the direction of the wind by arrows. The winds are predominantly westerly, as expected, and their spatial pattern represents a mid-latitude jet, whose axis is located at about 50°S between South America and Australia, and at about 55°S elsewhere; the strength of the jet in the latter region is somewhat weaker, with the exception of even weaker time-mean winds just east of South America. The standard deviation of the wind speed shown in panel (b) is fairly uniform throughout the Southern Ocean, with the most intense variability south of the stronger portion of the jet, and the weakest variance at the northern edge of the Southern Ocean. Modern data sets thus indicate that, at the $2^{\circ} \times 2^{\circ}$ resolution used here, the “furious fifties” are much more intense than the “roaring forties” of sailing days. Moreover, at this resolution and on a 5-yr average, winds off Cape Horn or the Cape of Good Hope are not particularly strong, although their standard deviation is maximal off Cape Horn and above the Agulhas Current, east and south of the Cape of Good Hope.

Color shading in panel (c) shows the distribution of the skewness of the zonal component of the surface wind, which

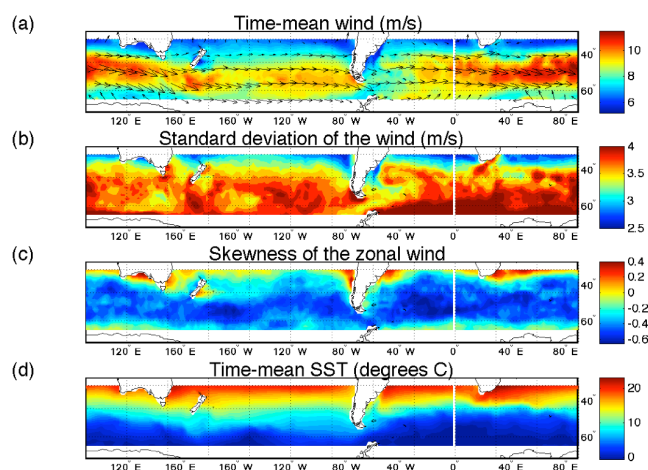


Fig. 2. Low-order moments of SLW and SST anomalies: (a) time-mean SLW; (b) standard deviation of SLW; (c) skewness of the zonal component of SLW; and (d) time-mean SST.

is found to be negative in the majority of the basin. Monahan (2004) explains this property of the zonal-wind anomalies in terms of a nonlinear surface drag law: according to this law, positive anomalies in u in the region with positive time-mean zonal winds will be subjected to stronger friction than negative anomalies, so that the resulting u -wind distribution will be negatively skewed.

The time-mean SST field (Fig. 2d) is consistent with the climatological wind (Fig. 2a) in that the strongest SST front is co-located with the strongest zonal jet, south of Africa and further eastward, at 40°S . This is presumably the region of the strongest ACC as well. The north-south SST gradients elsewhere are weaker. The SST variance (not shown) is largely uniform throughout the Southern ocean, with values around $1\text{--}2^{\circ}\text{C}$.

2.3.3 Principal component (PC) analysis

Prior to computing the EOFs and PCs of SLW and SST, we multiplied the time series of these quantities at each grid point by the square root of the cosine of its latitude, to account for the meridian convergence and get area-weighted grid-point contributions to the total variance of each field. The EOFs of SST and SLW were computed separately, and we used the combined (u, v) field to compute the latter. The percentages of variance accounted for by the first 100 EOFs of SLW are shown in Figs. 3a, b, while the analogous plots for SST EOF are in Figs. 3c, d.

Two leading SLW EOF pairs are somewhat separated from each other and from the rest of the modes (Fig. 3a) and together account for about 20 % of the total SLW variance (Fig. 3c). These two pairs are associated with the leading synoptic disturbances in the “weaker jet” ($160^{\circ}\text{W}\text{--}80^{\circ}\text{W}$) and “stronger jet” region ($40^{\circ}\text{W}\text{--}130^{\circ}\text{E}$), respectively; compare Fig. 2a and Fig. 4. Spatial analysis of Fig. 4a, b and 4c, d

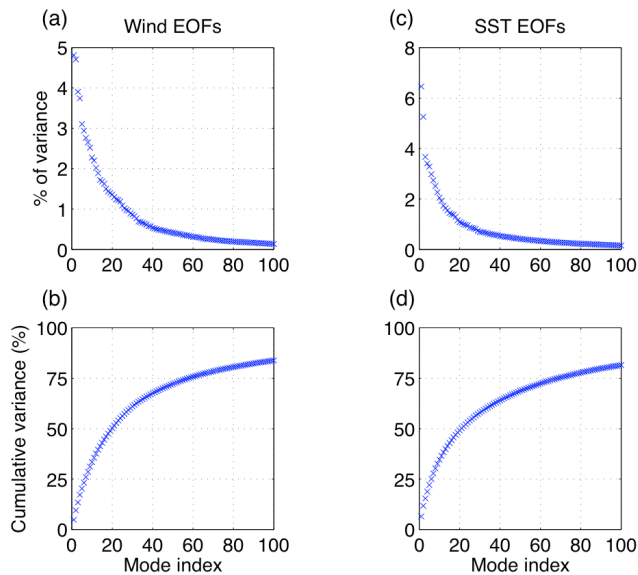


Fig. 3. Variances accounted for by the 100 leading EOF modes of (a, b) SLW and (c, d) SST; individual and cumulative variances appear in panels (a, c) and (b, d), respectively.

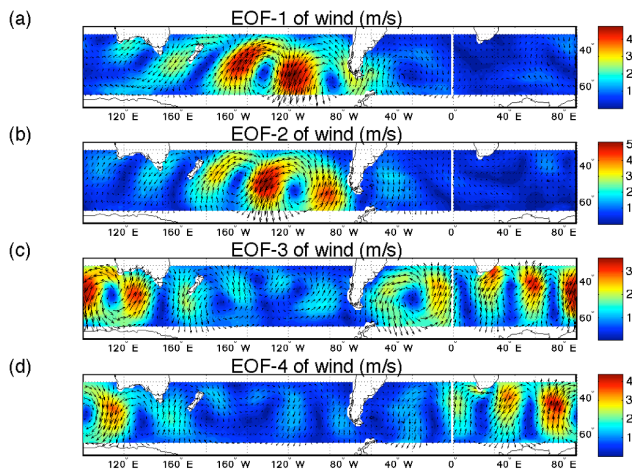


Fig. 4. Leading EOFs of SLW: (a, b) EOFs 1 and 2; (c, d) EOFs 3 and 4.

shows that these two modes are characterized by zonal wave numbers 8 and 9, respectively. For the SST, only two leading EOFs stand out from the rest (Fig. 3c), and account for about 12 % of the total SST variance (Fig. 3d). Both of these EOFs have a wavelike pattern with dominant zonal wavenumbers 3–4 (Figs. 5a, b) and pronounced interannual variability (Fig. 5c) suggesting their possible association with the ACW. In particular, if the time scale T of the ACW is set up by advection processes (Weisse et al., 1999), then $T = L/U$, where L and U are length and velocity scales, respectively. For wavenumber-3 patterns (Figs. 5a, b) $L \sim 6000$ km; given typical advective velocities of $U = 10 \text{ cm s}^{-1}$, one then ends up

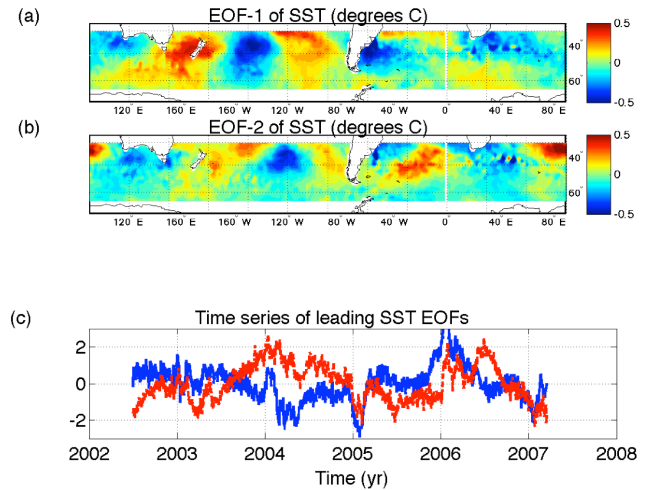


Fig. 5. Leading EOFs of SST: (a, b) EOFs 1 and 2; and (c) corresponding PCs (PCs 1 and 2 are shown as blue and red lines, respectively).

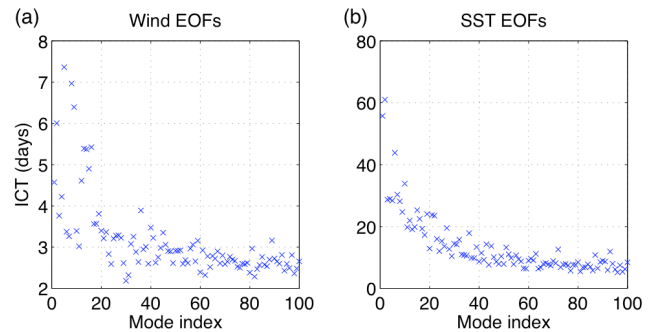


Fig. 6. Integral correlation time scales of the leading 100 PCs for: (a) wind, and (b) SST.

with the estimate $T \sim 2$ yr, consistent with Fig. 5c. The EOF spectrum becomes fairly flat roughly beyond mode 40 for SLW and mode 20 for SST. The leading 100 EOFs account for about 80 % of the total variance of both fields.

Figure 6 shows integral correlation time scales T_{int} of the leading 100 EOFs of SLW in panel (a) and SST in panel (b). The quantity T_{int} was defined as $T_{\text{int}} = \sum_{\tau=1}^{100} |c(\tau)| \Delta\tau$, where $c(\tau)$ is the autocorrelation of a given PC at the lag τ (in days), and $\Delta\tau = 1$ day. In general, the integral correlation time scale of the trailing modes is shorter than that of the leading modes, for both the SLW and SST PCs, although the dependence of T_{int} on the mode number is not monotonic. The leading EOF pairs of SLW have time scales of about 5.5 and 4 days, respectively, while the leading EOF pair of SST is characterized by $T_{\text{int}} \approx 60$ days. The latter estimate is an order of magnitude longer than the maximum T_{int} of the SLW EOFs, which is of about 7 days. Hasselmann (1976) introduced a null hypothesis for low-frequency variability SST anomalies, which involved integration of fast and essentially

random air-sea heat fluxes by an ocean mixed layer. The fast random heat flux forcing was associated with SLW variability, while the longer time scale of SST anomalies arose due to ocean mixed layer's thermal inertia. We argue that leading SST modes are not consistent with this null hypothesis for two reasons: (i) the Hasselmann mechanism is local, implying positive spatial correlations between SLW and SST patterns, whereas the patterns in Figs. 4 and 5a, b are not so correlated; and (ii) the interannual time scales of the leading SST modes (Fig. 5c) are longer than those associated with mixed layer thermal inertia. We thus conjecture that the leading SST modes arise from intrinsic ocean dynamics.

In fact, the arguments of the latter paragraph apply to most of the SST EOFs, more so for leading modes, and to a somewhat smaller degree for the trailing modes. The empirical stochastic models of SLW constructed in the next section will include the dependence on SST anomalies that span the subspace of their leading K EOFs, with $K = 50$ and $K = 75$. Therefore, any sensitivity of the SLW variability produced by empirical stochastic models to these SST anomalies should be interpreted as that caused by SST variability, rather than vice-versa.

3 Construction of empirical stochastic models

3.1 General methodology

We construct empirical stochastic models in the phase space of M leading EOFs of SLW, for various values of M (10–100), following the general methodology of Kravtsov et al. (2005b). In order to do so, we first form daily tendencies of N leading PCs of SLW: the tendency at day n , for example, is approximated as the difference between the value of a given PC at day $n+1$ minus the value of this PC at day n . The M time series of tendencies so obtained represent our response variables.

We will consider several versions of the empirical models. In the simplest, linear case, the main level of the empirical model is obtained by multiple linear regression (MLR) of each response variable onto M leading PCs of SLW, resulting in an equation of the form

$$\mathbf{x}^{n+1} - \mathbf{x}^n = \mathbf{B} \cdot \mathbf{x}^n + \mathbf{r}^n, \quad (1)$$

where \mathbf{x} is an M -component vector of the leading PCs, \mathbf{B} is an $M \times M$ matrix of the regression coefficients, while \mathbf{r} is the vector of M residual time series uncorrelated with each of the predictor variables; as before, n is the time index (in days). If we model \mathbf{r} as vector-noise $d\mathbf{w}$ that is white in time, but spatially correlated, then the formulation of Eq. (1) is a so-called linear inverse model (LIM) of SLW variability. The spatial correlation refers to that between the different components of the residual time series in \mathbf{r} (and $d\mathbf{w}$). Given random realizations of the forcing $d\mathbf{w}$, the LIM (1) can be integrated to produce surrogate time series of the M leading

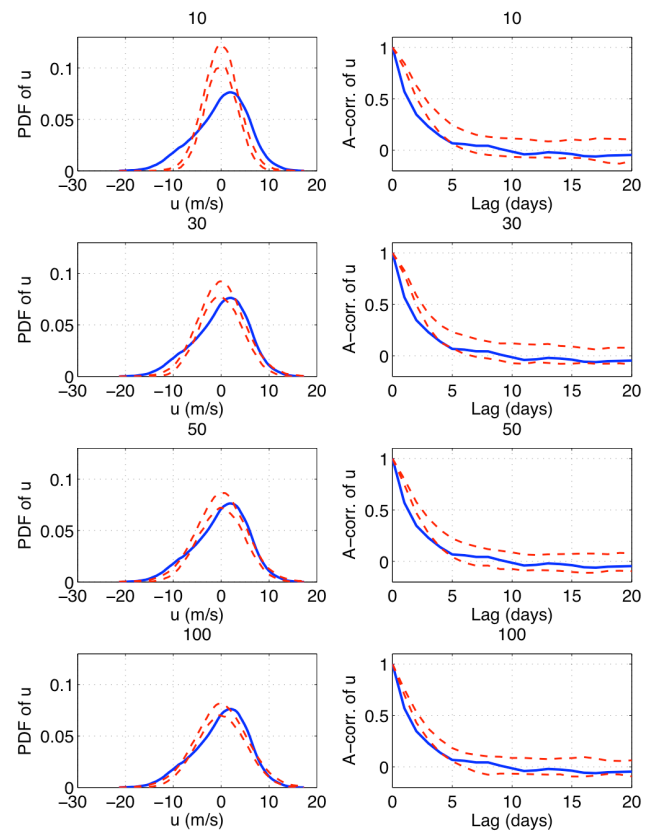


Fig. 7. Probability density function (PDF; left panels) and autocorrelation function (ACF; right panels) of the observed and simulated zonal velocity anomalies at 120° W and 55° S. Solid lines: the observed functions; dashed lines: 95 % spread based on SLW-only model with quadratic main level. The four rows show the results, from top to bottom, for the models constructed in the subspace of 10, 30, 50, and 100 PCs of SLW, respectively.

PCs of SLW, which can then be translated into variable SLW patterns in physical space by summing the SLW EOFs multiplied by the value of the corresponding surrogate PC at a given time. The statistics of such surrogate SLW realizations can then be compared to that of the observed anomalies to judge the performance of the LIM.

Kravtsov et al. (2005b) introduced several improvements to the LIM (1). In particular, it often happens that the autocorrelation of the residuals at nonzero lags is not negligible. In order to address this problem, Kravtsov et al. (2005b) proposed to construct an additional level of the inverse model; at this level, the tendencies of the main-level residuals are modeled as a linear function of the extended state vector $[\mathbf{x}^n; \mathbf{r}^n]$, consisting of the M original PCs, plus M first-level residuals:

$$\mathbf{r}^{n+1} - \mathbf{r}^n = \mathbf{B}_1 \cdot [\mathbf{x}^n; \mathbf{r}^n] + \mathbf{r}_1^n. \quad (2)$$

This exercise produces M second-level residuals: if the latter are not white in time, their tendencies can in turn be modeled as a linear function of the extended state vector consisting

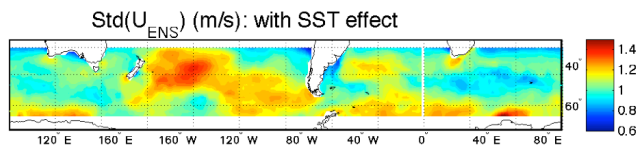


Fig. 8. The standard deviation of the wind speed time series obtained by taking the ensemble average of 100 simulations of SST-dependent SLW model forced by the observed history of SST anomalies. The model was constructed in the phase space of 100 leading EOFs of SLW. A typical (maximum) standard deviation of analogous SLW-only model's ensemble-mean time series (not shown) is 0.25 (0.55) – both values are smaller than the standard deviations shown here.

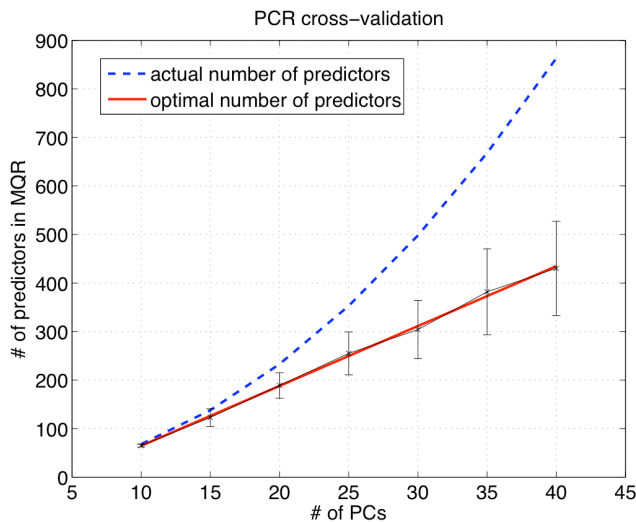


Fig. 9. The results of PCR cross-validation for the main level of our quadratic SLW-only models. The number on the abscissa shows how many SLW PCs are included in the model. The dashed line denotes the total number of predictors in the equation for each PC. The error bar plot (light solid line) shows the optimal number of PCR components, with the central value being the average of this number over its individual estimates obtained for each PC equation, and the bar representing the standard deviation of these estimates. The straight heavy line is the optimal linear fit of the dependence of the PCR-optimized number of components on the number of original variables (PCs) considered.

now of the M original PCs, M first-level residuals, as well as M second-level residuals. Additional levels can be added in the same way until the residual time series becomes white in time. Note that this procedure is different from merely modeling the main-level residual as colored noise, since it also takes into account any hidden dependency of the residual tendencies on the main-level PC predictors.

Another modification, which proved useful in modeling tropical SST evolution in Kondrashov et al. (2005), consisted of the inclusion of an explicit seasonal cycle. Despite our having removed the explicit seasonal cycle from the SLW

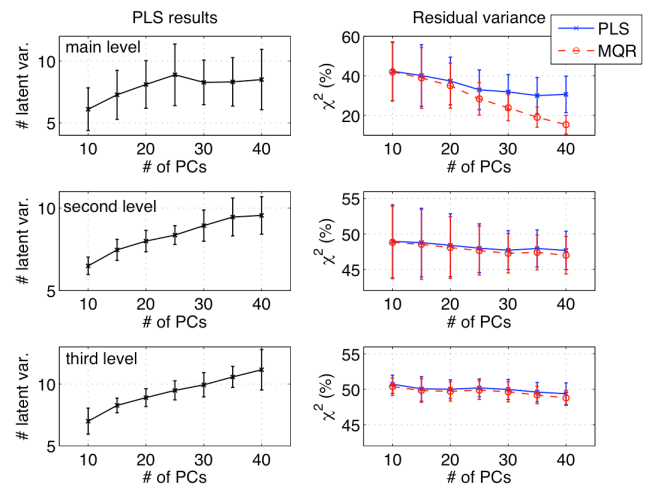


Fig. 10. PLS cross-validation results for the three-level empirical model with quadratic nonlinearity in the main level. The error bar plots show the mean and standard deviations of each quantity displayed computed using individual values of this quantity for each of the model equations (the number of equations is equal to the number of original PCs simulated by the empirical model). Left panels: the optimal number of PLS components; right panels: the percentage of variance unaccounted for by the regression; x -symbols show the results of PLS regression using the optimal number of latent variables, while the circles display the results of standard MQR, with all predictors considered.

fields prior to constructing our empirical stochastic model, the parameters of this model may still have some seasonal dependence. Kondrashov et al. (2005) found that the optimal way to incorporate such dependence is to include, at the main level of the model, two additional predictors, namely $\sin(2\pi t/365)$ and $\cos(2\pi t/365)$. The remainder of the procedure is unaltered, and the construction of the additional levels of the empirical stochastic models proceeds as described above.

Finally, the most significant modification of LIM methodology in Kravtsov et al. (2005b) was to consider nonlinear combinations of basic predictors. For example, one can include, in addition to M predictor variables (PC-1–PC- M), all possible quadratic combinations of PCs: the product of PC-1 with all of PC-1–PC- M , plus the product of PC-2 with PC-2–PC- M , and so on. Using index notations for vectors, matrices, and tensors, and assuming implicit summation over repeating indices, the modified main-level equation can be written as

$$x_i^{n+1} - x_i^n = a_{ijk} x_j^n x_k^n + b_{ij} x_j^n + c_i + r_i^n. \quad (3)$$

The coefficients of such a regression model are also found by the MLR procedure; however, since this procedure now employs an extended vector of predictor variables and their quadratic combinations, it is called multiple quadratic regression (MQR). Kravtsov et al. (2005b) argued that it is best to restrict the nonlinear modifications to the main level of the

empirical model, while the construction of additional levels proceeds as before. The main advantage of a nonlinear empirical model is that it can address nongaussian aspects of the observed variability. Its main disadvantage is a potentially much larger number of predictors: in a quadratic model based on M PCs and including two periodic seasonal cycle variables and a constant forcing term, the number of predictors is $M \times (M + 1)/2 + M + 3$, and so is the number of coefficients that need to be determined by the regression procedure for each of the M response variables. Kravtsov et al. (2005b) argued that this problem may be efficiently addressed by a variety of regularization procedures that allow one to avoid overfitting and construct nonlinear multi-level stochastic models with optimal predictive capabilities. We built on this approach here to develop a novel regularization algorithm for robust estimation of empirical model coefficients (see the appendices).

3.2 Stochastic model versions

Using the regularization methods described in appendix A, we have constructed several empirical model versions, which differed by the number of PCs considered, the order of nonlinearity at the main level of the model, and the presence or absence of the dependence on SST. All models had three levels, with levels 2 and 3 being linear. The SLW-only models with linear and quadratic main level were obtained in the subspace of $M = 10, 15, 20, 25, 30, 35, 40, 45, 50, 75$, and 100 PCs, and the cubic models for all of the above $M \leq 40$.

We found that the performance of all these model versions is very similar, for a given M . This result argues for using the linear SLW-only model, because it has the simplest form and the smallest number of coefficients; hence, it is more reliable and easily implemented than nonlinear models. The SLW model with SST dependence was based on $M = 100$ leading PCs of the SLW field – \mathbf{x} , and $L = 50$ or $L = 75$ leading PCs of the SST field – \mathbf{y} . The main level of the SST-dependent SLW model included linear dependence on SLW and SST PCs, the cross-product of each SLW PC with each SST PC, constant forcing term, and two seasonal cycle variables; no quadratic combinations of SLW PCs or SST PCs were used as predictors:

$$x_i^{n+1} - x_i^n = a_{ijk} x_j^n y_k^n + b_{ij}^x x_j^n + b_{ik}^y y_k^n + c_i^s \sin(2\pi t_n/365) + c_i^c \cos(2\pi t_n/365) + c_i + r_i^n. \quad (4)$$

4 Performance of empirical stochastic models

4.1 Simulation procedure

Empirical models constructed using the methodology described in the previous section and the appendices were used to produce 100 surrogate simulations of SLW variability, each of these simulations being 1719-day long. Despite the regularization applied when constructing regression models,

Table 1. Average number of initial-state resets for different empirical SLW-only models computed based on 100 surrogate simulations of a given model (see Sect. 4.1).

M	10	20	30	40	50	75	100
Linear	9	8	9	11	12	11	13
Quadratic	15	12	21	27	27	15	10
Cubic	9	25	35	48			

a few of the simulations using nonlinear models exhibited instability. In order to avoid such situations altogether, we have used the following procedure (Kravtsov et al., 2005b).

The models were integrated in ten-day chunks. The first chunk was started from random initial states. If at any time during this ten-day period the absolute value of any of the variables ended up outside their range (the latter ranges determined based on the values obtained for the training, model-construction period), then this ten-day simulation was discarded and restarted from another random state. The procedure was repeated as many times as necessary until a ten-day simulation with the values of all variables within the specified range was obtained. The final state from this simulation was then used to initialize the next ten-day simulation, for which the ranges of the variables were in turn monitored as before, and so on.

The threshold value for the PCs was computed as the observed maximum of the absolute value, over all the PCs and during the whole observational interval; the threshold values for the second and third-level variables were computed in the same way using “observed values” of these quantities. We kept track of the number of times the threshold condition above was violated, during each of 1719-day surrogate simulations. Table 1 lists the average values of this number for the simulations using linear, quadratic, and cubic SLW-only models. The average is computed over 100 available realizations of each empirical model.

In general, the number of initial-state resets is small, on the order of 10–20 resets during 1719-day-long simulation. Note that the resets do not always reflect instability – our linear model is, for example, always stable, in agreement with LIM theory (Penland, 1989, 1996; Penland and Ghil, 1993), but it still produces the values that exceed chosen thresholds from time to time. On the other hand, the cubic model for $M = 35$ and 40 does run out of control and may produce unbounded realizations of the simulated fields. Finally, the models that include SST forcing are not listed in Table 1, because they never produce realizations that exceed the threshold values. This fact suggests that coupling with SST is important for properly modeling SLW variability (see also Sect. 4.3).

4.2 Daily-to-monthly aspects of SLW variability

We illustrate the performance of the empirical models by examining first local aspects of the simulated SLW variability, based on the output of the quadratic, SLW-only model. Figure 7 shows probability density function (PDF; left panels) and autocorrelation function (ACF; right panels) of the observed and simulated zonal velocity anomalies at 120° W and 55° S – in the middle of an intense-jet region (see Fig. 2); the correspondence at other locations is qualitatively and quantitatively analogous. The heavy solid line in all the plots shows the observed PDF or ACF, while the dashed lines mark the 95 % spread in these quantities obtained from 100 realizations of the quadratic SLW model. The four top-to-bottom rows of Fig. 7 display the results from the empirical model based on $M = 10, 30, 50$, and 100 SLW PCs, respectively.

The empirical model of 10 leading PC components of SLW (upper row) produces a time series with a substantially smaller variance of the wind at the given location, while the time scale of SLW anomalies there is overestimated. Both of these results are to be expected, since the leading SLW EOFs account for a limited fraction of total variance (Fig. 3) and are generally characterized by the longest time scales (Fig. 6). Including progressively more components into the empirical model achieves continuous improvement of these two characteristics of SLW variability, with the 100-component model capturing quite well both the variance and the time scale of SLW anomalies. None of the model versions, however, captures the observed negative skewness of the zonal-velocity distribution. In fact, the quadratic model PDFs are essentially Gaussian and very similar to the ones obtained using simulations of the cubic and linear models (not shown).

We have tried a number of ways to better capture the skewness of the zonal-wind anomalies in our empirical models. These attempts included choosing a different EOF basis, which arranged the SLW patterns so that each of them would capture a significant fraction of variance, while having maximally skewed distribution, as well as blending our multi-level model methodology with the multiplicative-noise techniques of Sura and collaborators (see Sect. 1.2), but still failed to reproduce the negative skewness of the zonal-wind anomalies. We think that the reason for this failure is that the dynamics behind this negative skewness is essentially local, as it involves the effectively larger surface drag for positive u -wind anomalies in the region of the positive time-mean u -wind (Monahan, 2004). Considering the anomalies in the EOF basis does not optimally represent such local dynamics: Each of the PCs turns out to possess skewness values smaller than the typical skewness of the zonal wind at a certain grid point, and this skewness is identified by the regression procedure as negligible; hence, this non-Gaussian aspect of zonal-wind behavior is not properly represented in our empirical models. Non-local dynamics, though, are well represented in our statistical models of SLW evolution, as we will see in Sect. 4.3.

The correspondence between the observed and simulated statistics for meridional SLW components is similar or better than that in Fig. 7 (not shown), since the meridional wind distribution is generally more gaussian. Similar results are also obtained for other locations in the Southern Ocean (not shown). These local results are essentially indistinguishable between all versions of the empirical models including the SST-dependent version, given the number of SLW PCs considered.

4.3 SST effects on SLW evolution

We show here some preliminary evidence for the substantial oceanic imprint onto Southern Ocean's SLW variability; this oceanic effect is a necessary condition for the existence of active ocean-atmosphere coupling there. In order to do so, we have computed ensemble-averaged evolution of the SLW anomalies for a 100-member ensemble using the empirical stochastic model forced by the history of the observed SST anomalies, as well as this evolution for the SLW-only stochastic model. We then computed the standard deviation of the ensemble-averaged wind speed for both cases, at each grid point: the results of this computation for the SST-dependent SLW model are shown in Fig. 8. The standard deviation in the SST-dependent case is much larger (by a factor of 5–10), at all grid points, than that in the SLW-only case (not shown), and exhibits a distinctive large-scale spatial pattern, suggesting this SLW variability is forced by long-term, ocean-induced SST anomalies. We plan to address this intriguing behavior in a future paper (see Sect. 5).

In summary, the model constructed in the phase space of 100 leading EOFs of SLW and including, in addition, linear and bilinear interactions with SST anomalies restricted to the subspace of 75 leading EOFs of SST, as well as the seasonal effects, is stationary and captures several local and non-local aspects of SLW evolution, on all time scales. We plan to use this model as the atmospheric component of a hybrid coupled model in which the oceanic component will be a state-of-the-art GCM (see Sect. 5).

5 Summary and discussion

We have analyzed five years of remotely sensed data sets of sea-surface temperature (SST) and sea-level wind (SLW) over the Southern Ocean; the microwave sensors installed on recently launched NASA satellites provide an unprecedented quantity and quality of observations in the region. The missing data due to heavy rains or cloud coverage has been filled-in by singular spectrum analysis (SSA). The main technical outcome of this investigation is the construction of a statistical, stochastically forced model of SLW over the Southern Ocean; the model construction algorithm uses a number of essential innovations required to obtain robust estimates of the model's propagator. This model captures

detailed features of SLW variability on a wide range of time scales, from daily to interannual, and spatial scales spanning the range from the atmospheric Rossby radius to the basin scale. The model also accounts for ocean-atmosphere coupling via dependence of SLW equations on the SST anomalies.

The model's potential in helping to interpret observed evolution of Southern Ocean's climatic variables is briefly illustrated by identifying substantial oceanic imprint onto SLW variability, which may be indicative of possible coupled ocean-atmosphere effects in the Southern Ocean: ensemble averaging over 100 simulation of the statistical model forced by the observed SST anomalies reveals variability of a large magnitude and distinctive spatial pattern. The analogous ensemble average based on simulations of the SLW-only model is characterized by a very small magnitude and a lack of spatial coherence.

The construction of the above statistical models is rooted in the empirical methodology of Kravtsov et al. (2005b) and Kondrashov et al. (2005, 2006); however, the model construction algorithm is substantially modified and improved here in a number of ways that help choose the optimal model structure (see the appendices). These modifications make Kravtsov et al. (2005b) technique, previously used to identify low-dimensional behavior within high-dimensional noisy data, applicable to the analysis of the phenomena involving intermediate number of degrees of freedom. In particular, the most comprehensive statistical model operates in the subspace spanned by 100 leading empirical orthogonal functions (EOFs) of the daily SLW over the Southern Ocean, thus modeling the evolution of 100 corresponding principal components (PCs); the seasonal cycle was removed from all fields prior to performing the principal component analysis.

The model equations relate the time derivative of each PC to the right-hand side consisting of three parts: the part that depends on SLW only, the SST-dependent part, and the variable forcing term. The first part is approximated as a linear function of all PCs of the SLW field. The dependence on SSTs is modeled as the linear function of the leading 75 PCs of the SST, plus bilinear terms involving the cross-product of SLW and SST PCs; since this part is nonlinear, the seasonally dependent forcing term is also included. The variable forcing that drives the variability in the model is simulated in a separate set of equations that relate the time derivative of each component of the forcing vector to the linear function of SLW and SST PCs, as well as the forcing vector itself, and also include the second-level variable forcing. The second-level forcing's tendency is in turn modeled linearly in a way analogous to the main-level forcing, while the variable forcing at this last, third level of the model is approximated as spatially coherent noise that is white in time. The construction of this statistical model involved a novel multi-step regression algorithm to compute the coefficients of the model's propagator, as well as to determine the parameters of the noise.

Table 2. The number of statistically significant coefficients of a three-level quadratic inverse model based on M leading PCs of SLW (see Appendix B for further details).

Level	# of PCs (M)	# of all coeffs. (K)	# of significant coeffs. (K_s)	(K_s/K) × 100 %
Level 1	30	14 940	2248	15
	40	34 520	3426	10
	50	66 400	4333	7
	75	219 600	4849	2
	100	515 300	4660	1
Level 2	30	1800	695	39
	40	3200	994	31
	50	5000	1415	28
	75	11 250	2834	25
	100	20 000	4986	25
Level 3	30	2700	393	15
	40	4800	528	11
	50	7500	688	9
	75	16 875	1235	7
	100	30 000	1923	6

We plan to use the statistical model constructed in the present study to further investigate the dynamics of ocean-atmosphere interaction over the Southern Ocean. In particular, our current results may suggest the presence of active coupling in the region by identifying a nontrivial SLW response to the observed SST anomalies, although Bretherton and Battisti (2000) proposed alternative explanations to such findings. Goodman and Marshall (1999), on the other hand, formulated a theory of interannual-to-decadal coupled variability that is potentially applicable to the Southern Ocean. This theory predicts the existence of coupled modes, given a certain spatial phase relationship between SST patterns and SST-induced SLW anomalies; this phase relationship gives rise to Ekman pumping anomalies that force and modify the oceanic circulation and the associated SST field. It would be interesting to check whether we can detect such a phase relationship in our statistical model. Another very promising way to apply our empirical SLW model is to couple it to an oceanic GCM. We plan to achieve this coupling by blending the SST-dependent SLW model with atmospheric boundary layer model of Seager et al. (1995). The latter model needs the specification of boundary-layer winds to compute ocean-atmosphere heat fluxes. These winds will be supplied by the statistical model, and will also be used to compute the atmosphere-ocean momentum flux. The ocean model forced by heat, moisture, and momentum fluxes will predict the evolution of the SST field, which will, in turn, affect the future SLW anomalies. The experiments with such a hybrid coupled GCM of the Southern Ocean regions may provide invaluable insights into the dynamics of climate variability there.

Appendix A

PCR and PLS regression

The main regularization tool is cross-validation, in which one chooses randomly a subset of the vector time series (in the analyses below, we typically consider 80 % of original data points), applies a given regression technique, and then uses the regression model to reconstruct the segments of the time series that were omitted in the model identification step. The performance of the regression technique may then be assessed according, for example, to the smallness of the differences between the regression-based prediction and the actual values of the time series. We will use cross-validation in a number of different ways when constructing the empirical models below.

A major problem in applying MQR or MLR based on a large number of predictors is multi-collinearity (Press et al., 1994). This problem can be avoided by finding linear combinations of original predictors in such a way that their time series are uncorrelated, while each linear combination accounts for the maximum possible amount of the total variance. A natural way to determine this modified set of predictors is to apply principal component analysis to the original vector of predictors, and then use cross-validation for finding the optimal number of PCs to retain in the regression; this procedure is called the principal component regression (PCR). Note that since we construct our empirical models in the phase space of the data set's EOFs, the predictor variables in an LIM are already uncorrelated. On the other hand, the MQR predictors are the original set of PCs augmented by their quadratic combinations. Therefore, applying principal component analysis to this new multivariate data set generally produces a different set of predictors.

The PCR results for MQR based on several numbers of PCs, $M=10, 15, 20, 25, 30, 35$, and 40 , are displayed in Fig. 9; the values of M are shown on the abscissa of this graph. We computed the optimal number of PCR predictors for each of the M equations of the quadratic regression model. We thus obtained, for the model describing the evolution of M leading SLW PCs, M estimates of the optimal number of PCR components. The error bar plot in Fig. 9 shows the average value of this number over the M available estimates, along with its standard deviation. The dependence of the optimal number of PCR components on the number of original PCs is very well approximated by a linear fit (heavy solid line); this number is much smaller, for large M , than the maximum possible number of variables, which is equal, for MQR, to $M \times (M + 1)/2 + M + 3$.

PCR does a fairly good job in picking the smallest set of uncorrelated predictors that capture most of the variance. However, the choice of the PCR predictors does not involve at all the information about how well these predictors are correlated with the response variable. The procedure that does take into account this additional information is called partial

least-squares regression (PLS); see Abdi (2003) for a brief, but comprehensive review. We apply PLS to the set of optimal predictors determined via PCR cross-validation (Fig. 9), rather than to the original, much larger set of predictors.

Similarly to the PCR procedure, the leading PLS predictor is defined as a linear combination of the original predictor time series, but in this case the quantity being maximized is the correlation between this time series and the predictor time series. We found that applying PLS to each response variable individually produces better results than the matrix formulation of the PLS algorithm, which also considers linear combinations of all response variables and finds two sets of coefficients that define the mode of response and the mode of predictor variables that are maximally correlated (Abdi, 2003). In the general multivariate case, the weights of the leading PLS mode are found using Singular Value Decomposition (SVD) as the first right singular vector of the matrix $\mathbf{X}^T \mathbf{Y}$, where \mathbf{X} and \mathbf{Y} are the matrices whose columns are the time series of the predictor and response variables, respectively. The right singular vectors of $\mathbf{X}^T \mathbf{Y}$ define the weights for the response variables; in the univariate case, the single such weight is naturally equal to 1.

The time series of the leading PLS mode is obtained by summing the original time series of the predictor variables with the weights obtained as above. The signal associated with the leading PLS mode is then regressed out of both the response variable(s) time series, and all the predictor time series; this is done, once again, by only retaining the residual of the linear regression of each of these time series onto the time series associated with the leading PLS mode. The above procedure is then applied to the “reduced” response and predictor time series to obtain the next PLS mode, and so on to obtain all the PLS modes. The optimal number of modes to retain in this procedure is also determined by cross-validation.

The PLS cross-validation results for the main level of the quadratic models based on $M = 10, 15, 20, 25, 30, 35$, and 40 PCs are shown in the upper row of Fig. 10. The error bar plot in the left panel is analogous to that in Fig. 9, and shows, in this case, the optimal number of PLS components, which is found to be less than 10 for all M . The error bar plot with x -symbols (solid lines) in the right panel shows the residual variance as the percentage of the total response-variable variance; the expectation value and the standard deviation for a given M are, once again, based on the results of the PLS procedure applied to each of the M response variables (and, of course, the same set of original predictors). The additional error bar plot in the same panel (dashed line with circles) shows the same quantity based on the full MQR, which uses all of the original response variables. Note that for $M = 10, 15$, and 20 , only a few (definitely less than 10) effective predictor variables found by consecutive application of the PCR and PLS methodologies capture essentially the same amount of variance in the response variables as the MQR based on 63, 138, and 223 variables, respectively. For $M = 40$, the

residual variances differ by a factor of 2, which indicates that the additional variance “captured” by the original MQR procedure is associated with a substantial overfitting.

The additional panels in Fig. 10 show analogous results for the second ($2M$ original predictors) and third ($3M$ original predictors) level of the wind-only empirical stochastic model. The PCR pre-processing has not been applied to these levels, so that the PLS regularization acted directly on the original PCs and residuals. In each case, about a dozen optimal predictors are identified, which capture essentially the same amount of the response variance as the full MLR model for this level. Note that the residual variances become increasingly close to 50 % for the second and third level. Since our response variables have the form $r^{n+1}-r^n$ and the predictors include the term r^n , the case with no prediction skill (that is, r being pure white noise) will identify the regression coefficient multiplying r^n to be equal to -1 , and all other coefficients to be zero. In this case, the residual will be exactly equal to r^{n+1} , and therefore the residual variance will be exactly equal to the 50 % of the response-variable variance. The deviations of the residual variance from 50 % in the fourth level of the wind-only regression model are negligible (not shown), thus identifying the three-level empirical model to be optimal.

Appendix B

Selection of predictor variables

A few regression coefficients found by the application of PCR-and-PLS regularization, as described in appendix A, can be translated by trivial matrix manipulation into the coefficients of the empirical model in the original predictor-variable basis. Many of these coefficients are fairly small and do not contribute much to the predictive capability of a given empirical model. We therefore fine-tuned and enhanced our regression technique by the following procedure for the selection of the predictor variables.

This procedure was also based on subsampling of original predictor and response variables. For a model mimicking the evolution of M original PCs of SLW ($M = 10 - 100$), we first obtained 100 sets of regression coefficients by randomly applying PCR-and-PLS regularization to 100 randomly sampled subsets of the full original time series, each of which included 80 % of the original data points. The optimal number of PCR components in the quadratic model was estimated according to the linear approximation shown in Fig. 9. The general cubic model was also constructed for $M = 10 - 40$; for this model, we determined the optimal number of PCR components in a way analogous to that for the quadratic model, prior to applying the PLS regularization step. No PCR step was applied to the linear models. At the PLS step, we have used a fixed number of 25 latent variables to define the optimal subspace for regression. This number exceeded

the optimal one in Fig. 9 by at least a factor of two and thus could not result in underfitting. The regression coefficients so obtained were then translated into the original predictor-variable space.

If the interval between the 2nd and 97th percentile of a given regression coefficient obtained as described above contained the value zero, we excluded the corresponding predictor variable from consideration, thus forming a new, smaller subset of predictor variables. This subset was in turn subsampled 100 times and subjected to PCR-and-PLS regression to identify coefficients not significantly different from zero, and so on, until all coefficients of the final set of predictors were found to be significant. The same procedure was applied to the second and third level of each version of the inverse model. The final regression coefficients in each case were found by applying the PCR-and-PLS regularization to the fully sampled set of optimal predictors.

Table 2 lists the number of statistically significant nonzero coefficients of the three-level inverse model of M leading PCs of SLW; the main level includes quadratic nonlinearities and a seasonal cycle. The total number of coefficients at the main level is $(M \times (M + 1)/2 + M + 3) \times M$, at the second level $-2M^2$, and at the third level $-3M^2$. Note that the statistically significant coefficients are but a small fraction of the total number of coefficients. For example, for $M = 75$, the main level of the quadratic model has only 4849 nonzero coefficients, out of a maximum possible of 219600. This means that our regression procedure identified, on average, $4849/75 \approx 65$ nonzero coefficients in each of the 75 main-level equations; this number is an order of magnitude smaller than the number of degrees of freedom N_{DOF} in the time series of the length of 1719. If one estimates the decorrelation time scale of SLW anomalies to be 5 days, then $N_{\text{DOF}} \approx 1719/5 = 344 \gg 65$. Recall also, that the number of independent regression coefficients we have actually computed at each level is 25, which makes the number of coefficients/DOF comparison even more favorable.

Acknowledgements. This research was supported by NASA grant NNG-06AG66G-1 and DOE grant DE-FG02-02ER63413. IK was also supported by the NSF grant OCE-0749723 and SK – through University of Wisconsin-Milwaukee Research Growth Initiative program 2006–2007.

Edited by: A. Sterl

References

- Abdi, H.: Partial Least Squares (PLS) regression, In: Encyclopedia of Social Sciences Research Methods, edited by: Lewis-Beck, M., Bryman, A., and Futing, T., Sage, Thousand Oaks, CA, 2003.
- Arya, S. P.: Introduction to Micrometeorology, Academic Press, New York, 307 pp, 1988.
- Bretherton, C. S. and Battisti, D. S.: An interpretation of the results from atmospheric general circulation models forced by the time history of the observed sea surface temperature distribution, *Geophys. Res. Lett.*, 27, 767–770, 2000.
- Chelton, D. B., Esbensen, S. K., Schlax, M. G., Thum, N., Freilich, M. H., Wentz, F. J., Gentemann, C. L., McPhaden, M. J. and Schopf, P. S.: Observations of coupling between surface wind stress and sea surface temperature in the Eastern Tropical Pacific, *J. Climate*, 14, 1479–1498, 2001.
- Feldstein, S. B.: Is interannual zonal mean flow variability simply climate noise? *J. Climate*, 13, 2356–2362, 2000.
- Feliks, Y., Ghil, M., and Simonnet, E.: Low-frequency variability in the midlatitude atmosphere induced by an oceanic thermal front, *J. Atmos. Sci.*, 61, 961–981, 2004.
- Feliks, Y., Ghil, M., and Simonnet, E.: Low-frequency variability in the midlatitude baroclinic atmosphere induced by an oceanic thermal front, *J. Atmos. Sci.*, 64, 97–116, 2007.
- Garreaud, R. D. and Battisti, D. S.: Interannual (ENSO) and interdecadal (ENSO-like) variability in the Southern Hemisphere tropospheric circulation, *J. Climate*, 12, 2113–2123, 1999.
- Ghil, M., Allen, M. R., Dettinger, M. D., Ide, K., Kondrashov, D., Mann, M. E., Robertson, A. W., Saunders, A., Tian, Y., Varadi, F. and Yiou, P.: Advanced spectral methods for climatic time series, *Rev. Geophys.*, 40(1), 3.1–3.41, doi:10.1029/2000RG000092, 2002.
- Gill, A.: Atmosphere-Ocean Dynamics, Academic Press, 662 pp, 1982.
- Gille, S. T.: Float observations of the Southern Ocean. Part II: Eddy fluxes, *J. Phys. Oceanogr.*, 33, 1182–1196, 2003.
- Goodman, J. and Marshall, J.: A model of decadal middle-latitude atmosphere-ocean coupled modes, *J. Climate*, 12, 621–641, 1999.
- Grimm, A. M. and Silva Dias, P. L.: Analysis of tropical-extratropical interactions with influence functions of a barotropic model, *J. Atmos. Sci.*, 52, 3538–3555, 1995.
- Hall, A. and Visbeck, M.: Synchronous variability in the Southern Hemisphere atmosphere, sea ice, and ocean resulting from the annular mode, *J. Climate*, 15, 3043–3057, 2002.
- Hartmann, D. L.: A PV view of zonal flow vacillation, *J. Atmos. Sci.*, 52, 2561–2576, 1995.
- Hartmann, D. L. and Lo, F.: Wave-driven zonal flow vacillation in the Southern Hemisphere, *J. Atmos. Sci.*, 55, 1303–1315, 1998.
- Harrison, D. E. and Vecchi, G. A.: January 1999 Indian Ocean cooling event, *Geophys. Res. Lett.*, 28, 3717–3720, 2001.
- Hashizume, H., Xie, S.-P., Fujimara, M., Shiotani, M., Watanabe, T., Tanimoto, Y., Liu, W. T. and Takeuchi, K.: Direct observations of atmospheric boundary layer response to SST variations with tropical instability waves over the eastern Equatorial Pacific, *J. Climate*, 15, 3379–3393, 2002.
- Hasselmann, K.: Stochastic climate models. Part I: Theory, *Tellus*, 28, 473–485, 1976.
- Hsu, S. A.: Sea-breeze-like winds across the north wall of the Gulf Stream: An analytical model, *J. Geophys. Res.*, 89, 2025–2028, 1984.
- Jacobs, G. A. and Mitchell, J. L.: Ocean circulation variations associated with the Antarctic Circumpolar Wave, *Geophys. Res. Lett.*, 23, 2947–2950, 1996.
- Kamenkovich, I. V.: The role of daily surface forcing in setting the temperature and mixed layer structure of the Southern Ocean, *J. Geophys. Res.*, 101, C07006, doi:10.1029/2004JC002610, 2005.
- Karoly, D. J.: Southern hemisphere circulation features associated with El Niño-Southern Oscillation events, *J. Climate*, 2, 1239–1252, 1989.
- Karsten, R., Jones, H., and Marshall, J.: The role of eddy transfer in setting the stratification and transport of a circumpolar current, *J. Phys. Oceanogr.*, 32, 39–54, 2002.
- Kawanishi, T., Sezai, T., Ito, Y., Imaoka, K., Takeshima, T., Ishido, Y., Shibata, A., Miura, M., Inahata, H., and Spencer, R. W.: The Advanced Microwave Scanning Radiometer for the Earth Observing System (AMSR-E), NASDA's contribution to the EOS for Global Energy and Water Cycle Studies, *IEEE Trans Geosci. Remote Sens.*, 41, 184–194, 2003.
- Kondrashov, D. and Ghil, M.: Spatio-temporal filling of missing points in geophysical data sets, *Nonlin. Processes Geophys.*, 13, 151–159, 2006, <http://www.nonlin-processes-geophys.net/13/151/2006/>.
- Kondrashov, D., Kravtsov, S., and Ghil, M.: A hierarchy of data-based ENSO models, *J. Climate*, 18, 4425–4444, 2005.
- Kondrashov, D., Kravtsov, S., and Ghil, M.: Empirical mode reduction in a model of extratropical low-frequency variability, *J. Atmos. Sci.*, 63, 1859–1877, 2006.
- Koo, S. and Ghil, M.: Successive bifurcations in a simple model of atmospheric zonal-flow vacillation, *Chaos*, 12, 300–309, 2002.
- Koo, S., Robertson, A. W. and Ghil, M.: Multiple regimes and low-frequency oscillations in the Southern Hemisphere's zonal-mean flow, *J. Geophys. Res.*, 107, 4596, 13 pp., doi:10.1029/2001JD001353, 2002.
- Kravtsov, S., Robertson, A. W. and Ghil, M.: Bimodal behavior in the zonal mean flow of a baroclinic beta-channel model, *J. Atmos. Sci.*, 62, 1746–1769, doi:10.1175/JAS3443.1, 2005a.
- Kravtsov, S., Kondrashov, D., and Ghil, M.: Multi-level regression modeling of nonlinear processes: Derivation and applications to climatic variability, *J. Climate*, 18, 4404–4424, 2005b.
- Kravtsov, S., Robertson, A. W. and Ghil, M.: Multiple regimes and low-frequency oscillations in the Northern Hemisphere's zonal-mean flow, *J. Atmos. Sci.*, 63, 840–860, doi:10.1175/JAS3672.1, 2006a.
- Kravtsov, S., Berloff, P., Dewar, W. K., Ghil, M. and McWilliams, J. C.: Dynamical origin of low-frequency variability in a highly-nonlinear mid-latitude coupled model, *J. Climate*, 19, 6391–6408, doi:10.1175/JCLI3976.1, 2006b.
- Kummerow, C., Simpson, J., Thiele, O., Barnes, W., Chang, A. T. C., Stocker, E., Adler, R. F., Hou, A., Kakar, R., Wentz, F., Ashcroft, P., Kozu, T., Hong, Y., Okamoto, K., Iguchi, T., Kuroiwa, H., Im, E., Haddad, Z., Huffman, G., Ferrier, B., Olson, W. S., Zipser, E., Smith, E. A., Wilheit, T. T., North, G., Krishnamurti, T., and Nakamura, K.: The status of the Tropical Rainfall Measuring Mission (TRMM) after two years in orbit, *J. Appl. Meteorol.*, 39, 1965–1982, 2000.
- Kwok, R. and Comiso, J. C.: Southern Ocean climate and sea ice anomalies associated with the Southern Oscillation, *J. Climate*, 15, 487–501, 2002.

- Legras, B. and Ghil, M.: Persistent anomalies, blocking and variations in atmospheric predictability, *J. Atmos. Sci.*, 42, 433–471, 1985.
- Lindzen, R. S. and Nigham, S.: On the role of sea surface temperature gradients in forcing low-level winds and convergence in the tropics, *J. Atmos. Sci.*, 44, 2418–2436, 1987.
- Liu, W. T.: Progress on scatterometer application, *J. Oceanogr.*, 58, 121–136, 2002.
- Liu, W. T., Xie, X., Polito, P. S., Xie, S.-P. and Hashizume, H.: Atmospheric manifestation of tropical instability waves observed by QuikSCAT and Tropical Rain measuring Mission, *Geophys. Res. Lett.*, 27, 2545–2548, 2000.
- Lorenz, D. J. and Hartmann, D. L.: Eddy-zonal flow feedback in the Southern Hemisphere, *J. Atmos. Sci.*, 58, 3312–3327, 2001.
- Marshall, J. and Molteni, F.: Toward a dynamical understanding of atmospheric weather regimes, *J. Atmos. Sci.*, 50, 1792–1818, 1993.
- Meehl, G. A.: A reexamination of the mechanism of the semiannual oscillation in the Southern Hemisphere, *J. Climate*, 4, 911–926, 1991.
- Meehl, G. A., Hurrell, J. W. and van Loon, H.: A modulation of the mechanism of the semiannual oscillation in the Southern Hemisphere, *Tellus*, 50A, 442–450, 1998.
- Mitchell, T. P. and Wallace, J. M.: The annual cycle in equatorial convection and sea surface temperature, *J. Climate*, 5, 1140–1156, 1992.
- Mo, K. and Ghil, M.: Statistics and dynamics of persistent anomalies, *J. Atmos. Sci.*, 44, 877–901, 1987.
- Mo, K. and White, G. H.: Teleconnections in the Southern Hemisphere, *Mon. Wea. Rev.*, 113, 22–37, 1985.
- Monahan, A. H.: A simple model for skewness of global sea surface winds, *J. Atmos. Sci.*, 61, 2037–2049, 2004.
- Monahan, A. H.: The probability distribution of sea surface wind speeds. Part I: Theory and SeaWinds observations, *J. Climate*, 19, 497–520, 2006a.
- Monahan, A. H.: The probability distribution of sea surface wind speeds. Part II: Dataset intercomparison and seasonal variability, *J. Climate*, 19, 521–534, 2006b.
- O'Neill, L. W., Chelton, D. and Esbensen, S. K.: Observations of SST-induced perturbations of the wind stress field over the Southern Ocean on seasonal time scales, *J. Climate*, 16, 2340–2354, 2003.
- Penland, C.: Random forcing and forecasting using principal oscillation pattern analysis, *Mon. Wea. Rev.*, 117, 2165–2185, 1989.
- Penland, C.: A stochastic model of Indo-Pacific sea-surface temperature anomalies, *Physica D*, 98, 534–558, 1996.
- Penland, C. and Ghil, M.: Forecasting Northern Hemisphere 700-mb geopotential height anomalies using empirical normal modes, *Mon. Wea. Rev.*, 121, 2355–2372, 1993.
- Penland, C. and Sardeshmukh, P. D.: The optimal growth of tropical sea-surface temperature anomalies, *J. Climate*, 8, 1999–2024, 1995.
- Penland, C. and Matrosova, L.: Prediction of tropical Atlantic sea-surface temperatures using linear inverse modeling, *J. Climate*, 11, 483–496, 1998.
- Peterson, R. and White, W. B.: Slow oceanic teleconnections linking tropical ENSO and the Antarctic Circumpolar Wave, *J. Geophys. Res.*, 103, 24573–24583, 1998.
- Preisendorfer, R. W.: Principal Component Analysis in Meteorology and Oceanography, Elsevier, New York, pp. 425, 1988.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P.: Numerical Recipes, 2 Edn., Cambridge University Press, 994 pp., 1994.
- Reinhold, B. B. and Pierrehumbert, R. T.: Dynamics of weather regimes: Quasistationary waves and blocking, *Mon. Wea. Rev.*, 110, 1105–1145, 1982.
- Robinson, W.: A baroclinic mechanism for the eddy feedback on the zonal index, *J. Atmos. Sci.*, 57, 415–422, 2000.
- Seager R., Blumenthal, M. B. and Kushnir, Y.: An advective atmospheric mixed layer model for ocean modeling purposes: Global simulation of surface heat fluxes, *J. Climate*, 8, 1951–1964, 1995.
- Schmitz, W. J.: On the World Ocean Circulation: Volumes I, II. Woods Hole Oceanographic Institution Technical Report WHOI-96-08, 1996.
- Sura, P.: Stochastic analysis of Southern and Pacific Ocean sea surface winds, *J. Atmos. Sci.*, 60, 654–666, 2003.
- Sura, P. and Gille, S.: Interpreting wind-driven Southern Ocean variability in a stochastic framework, *J. Mar. Res.*, 61, 313–334, 2003.
- Sura, P. and Newman, M.: The impact of rapid wind variability upon air–sea thermal coupling, *J. Climate*, 21, 621–637, 2008.
- Sura, P. and Sardeshmukh, P. D.: A global view of non-gaussian SST variability, *J. Phys. Oceanogr.*, 38, 639–647, 2008.
- Sura, P., Newman, M., and Alexander, M. A.: Daily to decadal sea-surface temperature variability driven by state-dependent stochastic heat fluxes, *J. Phys. Oceanogr.*, 36, 1940–1958, 2006.
- Thompson, D. W. J. and Wallace, J. M.: Annular modes in the extratropical circulation. Part I: Month-to-month variability, *J. Climate*, 13, 1000–1016, 2000.
- Thompson, D. W. J., Wallace, J. M. and Hergerl, G. C.: Annular modes in the extratropical circulation, Part II: Trends, *J. Climate*, 13, 1018–1036, 2000.
- Van Loon, H.: The half-yearly oscillation in middle and high southern latitudes and the coreless winter, *J. Atmos. Sci.*, 24, 472–486, 1967.
- Van Loon, H.: Temperature, pressure, wind, cloudiness and precipitation in the Southern Hemisphere, In *Meteorology of the Southern Hemisphere*, Meteor. Monogr., No. 35, edited by: Newton, C. W., Amer. Meteor. Soc., 25–111, 1972.
- Vecchi, G. A. and Harrison, D. E.: Monsoon breaks and sub-seasonal sea surface temperature variability in the Bay of Bengal, *J. Climate*, 15, 1485–1493, 2002.
- Vecchi, G. A., Xie, S.-P. and Fischer, A. S.: Ocean–atmosphere covariability in the western Arabian Sea, *J. Climate*, 17, 1213–1224, 2003.
- Wallace, J. M., Mitchell, T. P., and Deser, C.: The influence of sea-surface temperature on surface wind in the Eastern Equatorial Pacific: Seasonal and interannual variability, *J. Climate*, 2, 1492–1499, 1989.
- Weisse, R., Mikolajewicz, U., Sterl, A. and Drijfout, S. S.: Stochastically forced variability in the Antarctic Circumpolar Current, *J. Geophys. Res.*, 104, 11049–11064, 1999.
- Wentz, F. J., Gentemann, C., Smith, D., and Chelton, D.: Satellite measurements of sea surface temperature through clouds, *Science*, 288, 847–850, 2000.
- Wetherill, G. B.: Regression Analysis with Applications, Chapman and Hall, 311 pp, 1986.

- Winkler, C. R., Newman, M., and Sardeshmukh, P. D.: A linear model of wintertime low-frequency variability. Part I: Formulation and forecast skill, *J. Climate*, 14, 4474–4494, 2001.
- White, W. B. and Peterson, R.: An Antarctic circumpolar wave in surface pressure, wind, temperature, and sea-ice extent, *Nature*, 380, 699–702, 1996.
- Wunsch, C.: Where do ocean eddy fluxes matter? *J. Geophys. Res.*, 104, 13235–13249, 1999.
- Yuan, X. and Martinson, D. G.: Antarctic sea ice extent variability and its global connectivity, *J. Climate*, 13, 1697–1717, 2000.
- Yuan, X. and Martinson, D. G.: The Antarctic dipole and its predictability, *Geophys. Res. Lett.*, 28, 3609–3612, 2001.