

On data-driven induction of the low-frequency variability in a coarse-resolution ocean model

E. A. Ryzhov^{a,c,*}, D. Kondrashov^{b,d}, N. Agarwal^a, J. C. McWilliams^b, and P. Berloff^a

^a*Department of Mathematics, Imperial College London, London, SW7 2AZ, UK*

^b*Department of Atmospheric and Oceanic Sciences, University of California, Los Angeles, CA 90095, USA*

^c*Pacific Oceanological Institute, Vladivostok, 690041, Russia*

^d*Institute of Applied Physics of the Russian Academy of Sciences, 603950, Nizhny Novgorod, Russia*

Abstract

This study makes progress towards a data-driven parameterization for mesoscale oceanic eddies. To demonstrate the concept and reveal accompanying caveats, we aimed at replacing a computationally expensive, standard high-resolution ocean model with its inexpensive low-resolution analogue augmented by the parameterization. We considered eddy-resolving and non-eddy-resolving double-gyre ocean circulation models characterized by drastically different solutions due to the nonlinear mesoscale eddy effects. The key step of the proposed approach is to extract from the high-resolution reference solution its eddy field varying in space and time, and then to use this information to improve the low-resolution analogue model.

By interactively coupling both the continuously supplied history of the eddy

*e.ryzhov@imperial.ac.uk

Email addresses: dkondras@atmos.ucla.edu (D. Kondrashov), n.agarwal17@imperial.ac.uk (N. Agarwal), jcm@atmos.ucla.edu (J. C. McWilliams), p.berloff@imperial.ac.uk (and P. Berloff)

field and the explicitly modelled low-resolution large-scale flow, we obtained the additional eddy forcing term which modified the low-resolution model and significantly augmented its solutions. This eddy forcing term represents the action of the eddy field, its coupling with the large-scale flow and is a key dynamical constraint imposed on the augmentation procedure.

Although the augmentation drastically improved the low-resolution circulation patterns, it did not recover the robust, intrinsic, large-scale low-frequency variability (LFV), which is an important feature of the high-resolution solution. This is by itself an important (negative) result that has significant implication for any data-driven eddy parameterization, especially, given the fact that we used the most complete information about the space-time history of the eddy fields. Note, when we supplied the reference (true) eddy forcing, rather than just the eddy field, the LFV was recovered. This suggests that the LFV is crucially dependent on the details of the space-time eddy forcing/large-scale flow correlations, which are not fully respected by the proposed augmentation procedure.

In order to overcome the deficiency and recover the LFV, we statistically filtered the augmented low-resolution model solution by projecting it onto the leading Empirical Orthogonal Functions (EOFs) of the large-scale component of the high-resolution reference solution. This operation allowed us to remove spurious effects associated with higher EOFs. We tested and confirmed that without using the data-driven eddy information this filtering alone cannot augment the low-resolution solution; but in conjunction with the eddy information, it produced desirable outcome.

Moreover, as a natural step towards parameterization, we took advantage of data-driven stochastic inverse modelling to obtain inexpensive emulators of the

eddy field and showed generally promising results of augmenting the coarse-resolution model with the obtained emulators. Our results showed that obtaining the LFV characteristics for the eddy parameterization, which is already capable of reproducing the large-scale flow pattern, should become a standard parameterization requirement, but it can be challenging to meet.

Keywords: Ocean dynamics, Mesoscale eddies, Eddy forcing, Parameterizations

1. Introduction

Numerical model solutions of complex oceanic flows are highly sensitive to the spatial grid resolution (Shevchenko and Berloff, 2015; Shevchenko et al., 2016). If the resolution is too coarse for representing mesoscale eddy dynamics, the resulting errors can be accumulated on large scales, which are nominally well-resolved even with dynamically coarse grids. On the one hand, this problem is now well understood in the ocean modeling community (Marshall et al., 2012; Bachman et al., 2017); on the other hand, resolving all the dynamically important scales is an insurmountable task, and many parameterizations aiming to circumvent this have been proposed and implemented (Gent and McWilliams, 1990; Frederiksen, 1999; Frederiksen et al., 2012; Porta Mana and Zanna, 2014; Berloff, 2015, 2016; Zanna et al., 2017; Berloff, 2018; Mak et al., 2018; Ryzhov et al., 2019). However, there is still no unified framework because different approaches are designed to account for different processes, and also each parameterization accounts for the effects of a certain range of scales.

Progress with parameterizations is hampered because the ocean circulation does not have spectral gaps between different ranges of scales; however, many theoretical insights rely on simple conceptual models with clear scale separation

19 (e.g., the Lorentz toy model (Majda et al., 1999; Fatkullin and Vanden-Eijnden,
20 2004; Kravtsov et al., 2005; Crommelin and Vanden-Eijnden, 2008; Arnold et al.,
21 2013; Chorin and Lu, 2015)). Furthermore, different scales are nonlinearly tangled
22 and accounting for this by understanding their interactions is difficult (Bachman
23 et al., 2017) but ultimately needed. The above-mentioned two aspects make the
24 problem of flow scale decomposition for the purposes of parameterizations open
25 and important. For now, the main constraint for a flow decomposition is rather
26 intuitive and vague: given the resolution of a coarse-grid model, we assume that
27 the unrepresented and dynamically distorted scales range from the Kolmogorov
28 scale to about 10 intervals of the computational grid; and the scales larger than the
29 grid interval are increasingly better accounted for by the model dynamics.

30 More specifically, in this paper we consider the classical, wind-driven, mid-
31 latitude ocean circulation model featuring two large-scale counter-rotating gyres
32 with the western boundary currents, and with their intense eastward jet exten-
33 sion that separates the gyres. Our focus is on the eastward jet region, where
34 the solutions of the model most critically depend on the spatial grid resolution
35 (Shevchenko and Berloff, 2015). With an inadequate resolution, misrepresenta-
36 tion of the mesoscale eddy dynamics results in an underdeveloped and even absent
37 eastward jet extension, whereas with a proper resolution, the eastward jet reappears
38 as a pronounced, meandering and vortex-shedding large-scale feature character-
39 ized by vigorous eddy dynamics and intensive eddy/large-scale interactions. Note,
40 that the flow decomposition into the large- and small-scale (i.e., mesoscale eddy)
41 components is not unique because of both the absence of the spectral gap and the
42 highly nonlinear dynamics — this complicates the analyses and parameterizations
43 of the eddy effects.

44 Our goal is to improve the analogue coarse-resolution double-gyre model by
45 feeding it with information obtained from solutions of the high-resolution model,
46 which is treated as the reference truth or the observed data. Ideally, this data-
47 driven approach should enable us to reproduce in the coarse-resolution model the
48 main characteristics of the high-resolution reference solution: (a) the large-scale
49 circulation pattern (specifically, the eastward jet extension with its adjacent recir-
50 culation zones) and (b) its intrinsic, large-scale low-frequency variability (LFV).
51 As we show in this paper, the latter characteristic proves more elusive to rectify,
52 even if the augmentation makes use of the full eddy information. To be precise,
53 one should aim at comparing the augmented coarse-resolution solution with the
54 large-scale component of the high-resolution solution, which is obtained by sta-
55 tistical filtering. Nevertheless, we focus on rectifying the large-scale circulation
56 patterns and LFV, which are interconnected, that are clearly transparent in the full
57 high-resolution solution as well, so we use it for the comparison.

58 Recently, Ryzhov et al. (2019) introduced a novel approach for augmenting
59 the coarse-resolution analogue model with data inferred from the high-resolution
60 truth; it involves the following main steps: (i) running the high-resolution model,
61 saving the solution data and verifying that the analogue low-resolution model
62 significantly misrepresents certain key features of the large-scale circulation; (ii)
63 decomposing the high-resolution data into some large-scale and small-scale (eddy)
64 fields; (iii) producing the eddy forcing term, which is based on the decomposed
65 fields and provides an important dynamical constraint, in order to exert extra
66 forcing and augment the low-resolution model in a dynamically consistent way.
67 Overall, an advantage of this approach is in combining its data-driven nature
68 with the transparent dynamical constraint, and this is strengthened by significant

69 flexibility of its practical implementations.

70 In this paper our goal is to extend the approach of (Ryzhov et al., 2019) by
71 significantly reducing and simplifying the information supplied from the high-
72 resolution reference truth. Now, instead of augmenting the model with the true
73 eddy forcing history coarse grained on the low-resolution grid, we supply only the
74 true eddy field (and its statistical emulation by a space-time stochastic process in a
75 separate experiment). This means that the eddy forcing term is now interactively
76 and continuously calculated *online* from the supplied eddy field history and the
77 dynamical low-resolution solution, which is treated as the prognostic large-scale
78 circulation. The approach is based on the implicit assumption that the low-
79 resolution model, if it is properly augmented, is adequate for representing the
80 large-scale circulation patterns and the LFV.

81 **2. Double-gyre model**

82 *2.1. Governing equations*

83 We use the same model configuration as in (Ryzhov et al., 2019). The model
84 has been extensively tested both in eddy-permitting and eddy-resolving regimes
85 (Marshall et al., 2012; Maddison et al., 2015; Shevchenko and Berloff, 2015;
86 Shevchenko et al., 2016; Ying et al., 2019). A brief description is as follows. The
87 quasi-geostrophic (QG) potential vorticity (PV) evolution in 3 stacked isopycnal
88 layers ($i = 1..3$ from top to bottom) with densities ρ_i ($\rho_1 = 1000$, $\rho_2 = 1001.498$,
89 $\rho_3 = 1001.62$ kg m⁻³) and heights H_i ($H_1 = 250$, $H_2 = 750$, $H_3 = 3000$ m) is
90 given by

$$\frac{\partial q_i}{\partial t} + J(\psi_i, q_i) + \beta \frac{\partial \psi_i}{\partial x} = \frac{W(x, y)}{\rho_i H_i} \delta_{1i} - \gamma \Delta \psi_i \delta_{3i} + \nu \Delta^2 \psi_i, \quad (1)$$

91 where q_i is the PV anomaly, ψ_i is the streamfunction, $J(\cdot, \cdot)$ is the Jacobian operator,
 92 δ_{ij} is the Kronecker delta, Δ is the horizontal Laplacian, $\beta = 2 \cdot 10^{-11} \text{ m}^{-1} \text{ s}^{-1}$ is
 93 the planetary vorticity gradient, ν is the eddy viscosity (varies for different spatial
 94 resolutions used in the study), $\gamma = 4 \cdot 10^{-8} \text{ s}^{-1}$ is the bottom friction parameter.
 95 The basin is north-south oriented square $-L \leq x, y \leq L$, where $2L = 3840 \text{ km}$.

96 The upper-ocean layer is forced by the stationary asymmetric wind stress curl

$$W(x, y) = \begin{cases} -\frac{\pi\tau_0 A}{L} \sin \frac{\pi(L+y)}{L+Bx}, & y \leq Bx, \\ \frac{\pi\tau_0}{LA} \sin \frac{\pi(y-Bx)}{L-Bx}, & y > Bx, \end{cases} \quad (2)$$

97 where the asymmetry, tilt, and wind stress magnitude parameters are $A = 0.9$,
 98 $B = 0.2$, and $\tau_0 = 0.08 \text{ N m}^{-2}$, respectively.

99 The PV anomalies and streamfunctions are related through

$$\begin{aligned} q_1 &= \Delta\psi_1 + S_1(\psi_2 - \psi_1), \\ q_2 &= \Delta\psi_2 + S_{21}(\psi_1 - \psi_2) + S_{22}(\psi_3 - \psi_2), \\ q_3 &= \Delta\psi_3 + S_3(\psi_2 - \psi_3), \end{aligned} \quad (3)$$

100 where the stratification parameters S_1, S_{21}, S_{22}, S_3 are chosen to yield the first and
 101 second baroclinic Rossby deformation radii of 40 and 23 km, respectively. The
 102 boundary conditions are no-flow-through and partial-slip (with the partial-slip
 103 length scale equal to 120 km); the mass is conserved in each layer. The model is
 104 solved using the high-resolution CABARET method that features a second-order,
 105 non-dissipative and low-dispersive, conservative advection scheme (Karabasov
 106 et al., 2009).

107 Given an adequately fine spatial resolution, the model is capable of resolving
 108 the eddies that maintain the well-developed eastward jet extension of the western
 109 boundary current. Otherwise, the eastward jet extension is under-predicted or even

110 absent because the backscatter process of the energy transfer from the eddies to
111 the large-scale flow is under-resolved by the model (Jansen and Held, 2014; Jansen
112 et al., 2015; Shevchenko and Berloff, 2016; Berloff, 2018).

113 2.2. *Differences of flow structures in eddy-resolving and eddy-permitting regimes*

114 We consider two spatial grid resolutions for simulating the eddy-permitting
115 (low-resolution) and eddy-resolving (high-resolution) flow regimes: 129×129
116 and 513×513 , respectively. For resolving the western boundary layer (Berloff
117 and McWilliams, 1999), the low-resolution configuration is run with the viscosity
118 $\nu = 50 \text{ m}^2 \text{ s}^{-1}$, whilst the high-resolution one has $\nu = 2 \text{ m}^2 \text{ s}^{-1}$. In both cases,
119 the model is first spun-up for 100 years until a statistically equilibrated state is
120 achieved; then, its daily output is saved for another 90 years for further analyses.

121 The differences in the resulting flows are well-documented (Shevchenko and
122 Berloff, 2015; Ryzhov et al., 2019), so here we only note that the low-resolution
123 model does not induce a proper eastward jet extension (Fig. 1a), whereas the
124 high-resolution one features a well-pronounced, eddy-driven eastward jet with
125 the adjacent recirculation zones (Fig. 1b). Throughout the paper we make use
126 of the standard deviation instead of the time-mean when address the problem of
127 rectifying the large-scale circulation patterns. The standard deviation accentuates
128 more saliently the differences also easily seen in time-mean patterns.

129 Not only the spatial patterns but also the temporal variabilities of the reference
130 solutions are different. To reveal details of the latter, we used the Data-Adaptive
131 Harmonic Decomposition (DAHD) method (Chekroun and Kondrashov, 2017;
132 Kondrashov et al., 2018), which characterizes a complex and multiscale spatio-
133 temporal variability by extracting spatial data-adaptive harmonic modes (DAHMs)
134 such that each one of them oscillates at a single temporal frequency and is spatially

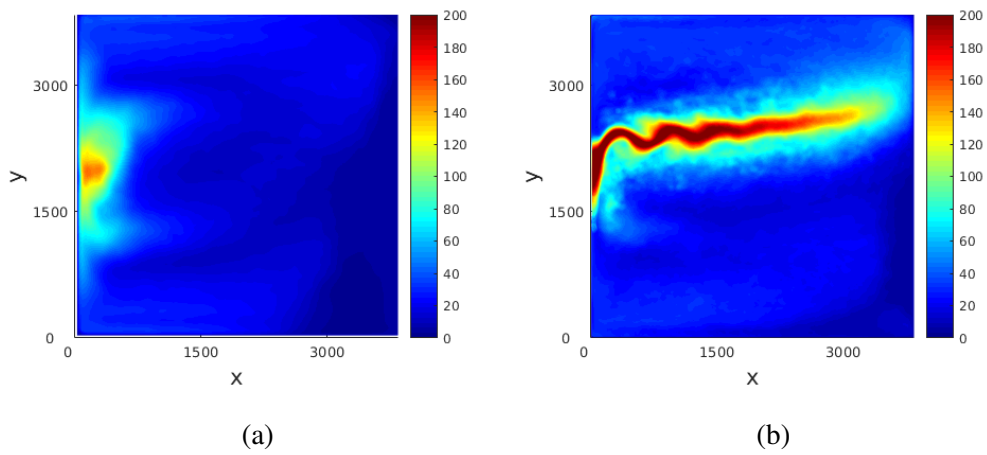


Figure 1: Standard deviation of the upper-layer PV anomaly (q_1) produced by the (a) low-resolution (129^2) and (b) high-resolution (513^2) models. The solutions emphasise the crucial effect of the spatial resolution. Nondimensional color scale units (PV is normalized using the length scale 3×10^4 m, corresponding to the low-resolution grid interval, and the velocity scale 0.01 m/s) are the same across all the figures.

135 orthogonal to all other modes at that frequency (see Appendix A for details).
 136 The DAHD has been successfully applied to characterize variabilities in different
 137 geophysical datasets including ocean circulation (Kondrashov et al., 2018; Ryzhov
 138 et al., 2019; Kondrashov et al., 2020), sea ice (Kondrashov et al., 2018a,b), and
 139 space physics (Kondrashov and Chekroun, 2018).

140 Here, we applied the DAHD to the upper-ocean PV anomaly fields of the
 141 reference solutions. To make our analysis computationally tractable, first, these
 142 fields were compressed using the standard principal component analysis (PCA)
 143 (Preisendorfer, 1988) to retain the leading $d = 2000$ empirical orthogonal function
 144 (EOF) modes. These modes capture 98% and 95% of the variance in the low- and
 145 high-resolution solutions, respectively. Next, the original PV anomaly fields were
 146 projected onto the retained EOFs to obtain the corresponding principal components
 147 (PCs). These $d = 2000$ PCs were used as inputs for the DAHD frequency-domain
 148 formulation, which is tailored for analysis of high-dimensional datasets (Chekroun
 149 and Kondrashov, 2017; Ryzhov et al., 2019) and based on the singular value
 150 decomposition (SVD) of the $d \times d$ symmetrized complex cross-spectral matrix
 151 $\mathfrak{S}(f)$:

$$\mathfrak{S}_{p,q} = \begin{cases} \widehat{\rho^{p,q}}(f) & \text{if } q \geq p, \\ \widehat{\rho^{q,p}}(f) & \text{if } q < p, \end{cases} \quad (4)$$

152 where $1 \leq p, q \leq d$; and $\widehat{\rho^{p,q}}(f)$ is the Fourier transform of the double-sided cross-
 153 correlation coefficients $\rho^{(p,q)}(m)$ estimated for all pairs of the channels (PCs) p
 154 and q , and for the time lag m , up to its maximum $M - 1$; i.e. $-(M - 1) \leq$
 155 $m \leq M - 1$. Each singular value $\sigma_k(f)$ of $\mathfrak{S}(f)$ is associated with a pair
 156 of negative/positive eigenvalues $(\lambda_k^+(f), \lambda_k^-(f))$ obtained by using the standard
 157 DAHD time-domain formulation and an eigen-decomposition of a matrix formed of

158 the elements $\rho^{(p,q)}(m)$ (Kondrashov et al. (2018); Ryzhov et al. (2019); Kondrashov
 159 et al. (2020)):

$$\lambda_k^+(f) = -\lambda_k^-(f) = \sigma_k(f), \quad 1 \leq k \leq d, \quad (5)$$

160 The DAHD power spectrum is obtained by plotting eigenvalues $|\lambda(f)|$ which
 161 represent energy conveyed by associated DAHMs; the frequency f is equally
 162 spaced with the Nyquist interval $[0, 0.5]$ across the M values:

$$f = 0.5 \frac{(\ell - 1)}{M - 1}, \quad \ell = 1, \dots, M. \quad (6)$$

163 The adequate spectral resolution in the low-frequency part is achieved by
 164 considering 30K days long PCs, sub-sampled every 5 days. Thus, we have $N =$
 165 6000 samples and use the largest possible embedding window $M = N/2 = 3000$
 166 for the maximum spectral resolution in the frequency domain.

167 Despite the overall similarity of the DAHD spectra shown in Fig. 2 and char-
 168 acterized by the bands of higher values separated by the gaps from the broadly
 169 distributed bands of lower values, as well as by the power-law behaviors in the high-
 170 frequency range, the low-resolution solution spectrum has significantly smaller
 171 magnitudes, which indicate the reduced eddy activity. In the upper band, there
 172 are two $|\lambda|$ values at each frequency, each of them corresponding to a negative-
 173 positive pair (see Eq.5). The observed gap in the spectrum can be interpreted as a
 174 dominance of a particular physical mechanism of energy distribution and transfer
 175 across all the temporal frequencies. However, the exact interpretation of the spec-
 176 tra is significantly hindered by the nonlinear character of the underlying physical
 177 interactions. Here, we use the spectra to diagnose the LFV and its profound effect
 178 on the spectrum.

179 The striking difference is the pronounced LFV in the high-resolution solution
 180 (see the blue dots in Fig. 2b at the period ≈ 17 years), and its complete absence

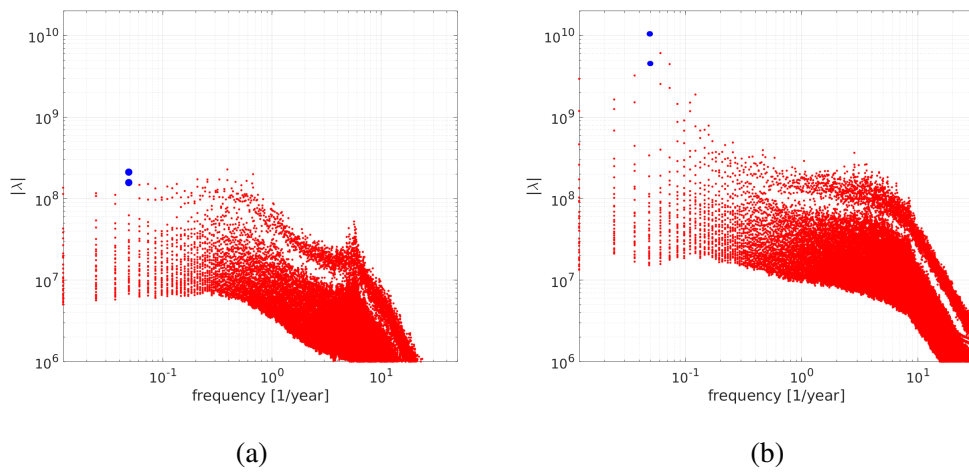


Figure 2: Temporal spectral content of the reference solutions with: (a) 129^2 and (b) 513^2 grids. Shown are the 30 largest values of $|\lambda|$ per frequency, as given by the DAHD power spectrum of the upper-layer PV anomalies. The blue dots in panel (b) indicate maximum of the broadband spectral peak corresponding to the low-frequency variability (LFV) $\approx 17\text{yr}$ in the high-resolution solution; this LFV is absent in the low-resolution solution (panel (a)).

181 in the low-resolution solution (Fig. 2a). This interdecadal LFV was studied
182 elsewhere (Berloff and McWilliams, 1999; Berloff et al., 2007; Shevchenko et al.,
183 2016), and here we just note that the quality of an augmented low-resolution model
184 can be tested by the model's capability to simulate this LFV.

185 *2.3. Low-frequency variability as an indicator of properly resolved small scales*

186 As we pointed out in the previous section, one of the most remarkable dynamical
187 features which differentiate the low- and high-resolution solutions is the LFV in
188 the latter. The LFV manifests itself as the total energy modulation with the period
189 ≈ 17 years (Berloff and McWilliams, 1999; Kondrashov and Berloff, 2015). A
190 peculiar characteristic of the LFV is that it appears only if the double-gyre model
191 resolves the eddies and hence activates the essential eddy backscatter mechanism
192 (Berloff et al., 2007; Shevchenko and Berloff, 2016). The backscatter here means
193 that the energy from the small scales is transferred to the large scales and thus
194 impacts the large-scale circulation. If the spatial resolution is too coarse (even in
195 eddy-permitting regimes), the small scales are not resolved and in turn the large
196 scales are also under-saturated, which introduces many inconsistencies in the flow
197 when comparing solutions corresponding to differing spatial resolutions.

198 Ryzhov et al. (2019) demonstrated that the low-resolution model is in principle
199 capable of inducing the LFV, provided that it is augmented with the eddy forcing
200 history provided by the high-resolution data. Our goal now is to reduce the amount
201 of the information inferred from the high-resolution data, but still be able to capture
202 the LFV and induce it in the augmented low-resolution model. Thus, instead of
203 using the complete high-resolution data for estimating the true eddy forcing and
204 using it to augment the low-resolution model, we intend to use only the true eddy
205 component of the flow, and to calculate the augmenting eddy forcing interactively

206 by using the large-scale flow predicted by the augmented low-resolution model.

207 **3. Scale decomposition of the high-resolution solution**

208 The high-resolution solution, which is treated as the truth, should be decom-
209 posed into a combination of large-scale and small-scale (eddy) components. The
210 former one should be adequately captured by an augmented low-resolution model;
211 whilst the latter one may remain largely unresolved. However, we know that the
212 true eddy forcing adequately augments the low-resolution model, and this is a
213 necessary condition for our next steps.

214 An issue of significant concern is that the large-scale/eddy flow decomposi-
215 tion, which is central to the proposed augmentation scenarios, is neither unique
216 nor clearly constrained by dynamical or statistical arguments. For now, various
217 methods assume (Hasselmann, 1988; von Storch et al., 1995; Schmid, 2010; Li
218 and von Storch, 2013; Dijkstra, 2013, 2018; Viebahn et al., 2019; Agarwal et al.,
219 2020) that the implemented flow decomposition (i.e., scale separation) is practi-
220 cally meaningful, and then build upon this assumption; our work is fully within
221 this framework.

222 A formal scale decomposition for an arbitrary 2D time-dependent field Ξ (in
223 our case, Ξ stands for the layer-wise streamfunctions ψ_i and PV anomalies q_i)
224 reads

$$\Xi(x, y, t) = \bar{\Xi}(x, y, t) + \Xi'(x, y, t), \quad (7)$$

225 where the overbar and prime indicate the large-scale and eddy components, respec-
226 tively. With this in mind, we decomposed the high-resolution streamfunctions ψ_i
227 by the moving-average square filter of size W ; and the corresponding PV anoma-
228 lies are obtained by differentiation (akin eq. (3)). We justify our choice of W

229 by focusing on mesoscale eddies, which are scaled by the first baroclinic Rossby
 230 deformation radius, but we also admit that the problem contains many length scales
 231 and they vary geographically making the flow decomposition a difficult and open
 232 problem. The problem stems from the fact that for linear flows (when all the active
 233 scales are well separated in the Fourier spectra), the filter size should linearly de-
 234 pend on the ratio between the fine - and coarse - resolution grids. However, in our
 235 case, there is no separation between the active scales and the filter size is chosen
 236 based on the expected dynamical features we would like to filter out assuming the
 237 coarse-resolution model being unable to resolve them. In our case, these features
 238 are mesoscale eddies with length scales of order of the first baroclinic Rossby
 239 deformation radius ($\approx 10 - 100$ km).

240 Preliminary analyses (Ryzhov et al., 2019) suggest that the filter size of $W = 21$
 241 of high-resolution grid intervals (≈ 150 km in physical units) is adequate, but we
 242 also tested $W = 41$ as a tribute to the unavoidable sensitivity analysis. The
 243 eddy fields (calculated on the high-resolution spatial grid 513×513) were coarse-
 244 grained to be fed into the low-resolution (129×129) model by averaging over four
 245 adjacent grid cells in each spatial direction.

246 Guided by the fact that the LFV is eddy-driven, we substituted (7) into the
 247 governing equation (1) and for each layer obtained:

$$\frac{\partial \bar{q}_i}{\partial t} + J(\bar{\psi}_i, \bar{q}_i) = \mathcal{F}_i(\bar{\psi}_i, \bar{q}_i, \psi'_i, q'_i) + \mathcal{H}_i(\bar{\psi}_i, \bar{q}_i) + \mathcal{L}_i(\psi'_i, q'_i), \quad (8)$$

248 where the operator \mathcal{H}_i contains all terms involving only the large-scale components;
 249 the linear operator \mathcal{L}_i contains the eddy tendency term and all linear terms involving
 250 the eddy components; and the remaining term,

$$\mathcal{F}_i = - (J(\bar{\psi}_i, q'_i) + J(\psi'_i, \bar{q}_i) + J(\psi'_i, q'_i)) , \quad (9)$$

251 is the eddy forcing (Berloff, 2005) due to nonlinear coupling of the large-scale and
 252 eddy components. The linear eddy term \mathcal{L}_i can be neglected, since its contribution
 253 to the eastward jet (as we checked) is about 2% of that of the eddy forcing.

254 Ryzhov et al. (2019) established that the eddy-forcing term, when properly
 255 preprocessed with respect to the low-resolution dynamics, can be effectively added
 256 into the low-resolution model to improve significantly the mean flow and transient
 257 (spectrally treated) characteristics of its solutions. In this work, our goal is to
 258 reduce the amount of the high-resolution information by feeding the eddies rather
 259 than the eddy forcing information (which depends on both the eddies and large
 260 scales) into the augmented model.

261 **4. Feeding the eddy field into the low-resolution model**

262 With only the eddies being fed to the augmented model, the external infor-
 263 mation is subtler, which makes it harder for the low-resolution model to resolve
 264 desired dynamics resembling the fine-resolution reference solution such that the
 265 eastward extension of the jet is noticeably rectified and the low-frequency variabil-
 266 ity is present. At the same time, gauging the possibility of reducing the amount
 267 of data necessary for successful parameterization and errors introduced due to the
 268 incompleteness of the data is practically important.

269 The governing equations for the augmented low-resolution model are, thus:

$$\frac{\partial q_i}{\partial t} + J(\psi_i, q_i) = \mathcal{F}_i(\psi_i, q_i, \psi'_i, q'_i) + \mathcal{H}_i(\psi_i, q_i), \quad (10)$$

270 where the small-scale (eddy) fields ψ'_i, q'_i are taken from the high-resolution data,
 271 and the prognostic low-resolution, large-scale variables ψ_i, q_i are continuously
 272 updated *online* during numerical integration of the model. We used all 90 years

273 of the daily output to extract the eddy fields and then linearly interpolated them in
274 time in-between the data records. An important issue of determining the minimal
275 length of the eddy history for the quality augmentation of the low-resolution model
276 is left outside the scope of the paper and will be addressed elsewhere.

277 We assessed the quality of the augmented low-resolution solution by looking
278 into the simulated eastward jet region, focusing on its large-scale circulation pat-
279 terns (evinced by the standard deviation in time) and LFV. The augmented-model
280 eastward jet has improved but is still substantially different from the reference truth,
281 as can be seen by comparing Fig. 3a and Fig. 1a. Similarly large discrepancies are
282 seen in the augmented-model DAHD spectrum (Fig. 3b), which completely lacks
283 the LFV. The interactive eddy forcing (Fig. 4a) can be significantly less efficient
284 because it is noticeably weaker than the true eddy forcing (Fig. 4c). We checked
285 this by considering the more energetic eddy field extracted with the larger filter size
286 $W = 41$ (Fig. 4b), but although the resulting eddy forcing is as intensive as the
287 true one, the augmented model is still incapable of generating the LFV as implied
288 by the DAHD spectrum (Fig. 3d). From this, we conclude that feeding even the
289 most complete eddy fields into the model is still not sufficient for augmenting the
290 solution. So, one has to use additional information from the high-resolution data
291 to induce the LFV.

292 It has been already established (Ryzhov et al., 2019) that the true (off-line)
293 eddy-forcing (Fig. 4b) generates the LFV in the augmented solution; therefore,
294 we know that one way or another the model can be successfully augmented with
295 the right amount of the extra information. One way to add this information is
296 by interactively projecting the augmented solution onto the leading, true large-
297 scale EOFs, and this can be viewed as a weak statistical constraint imposed by

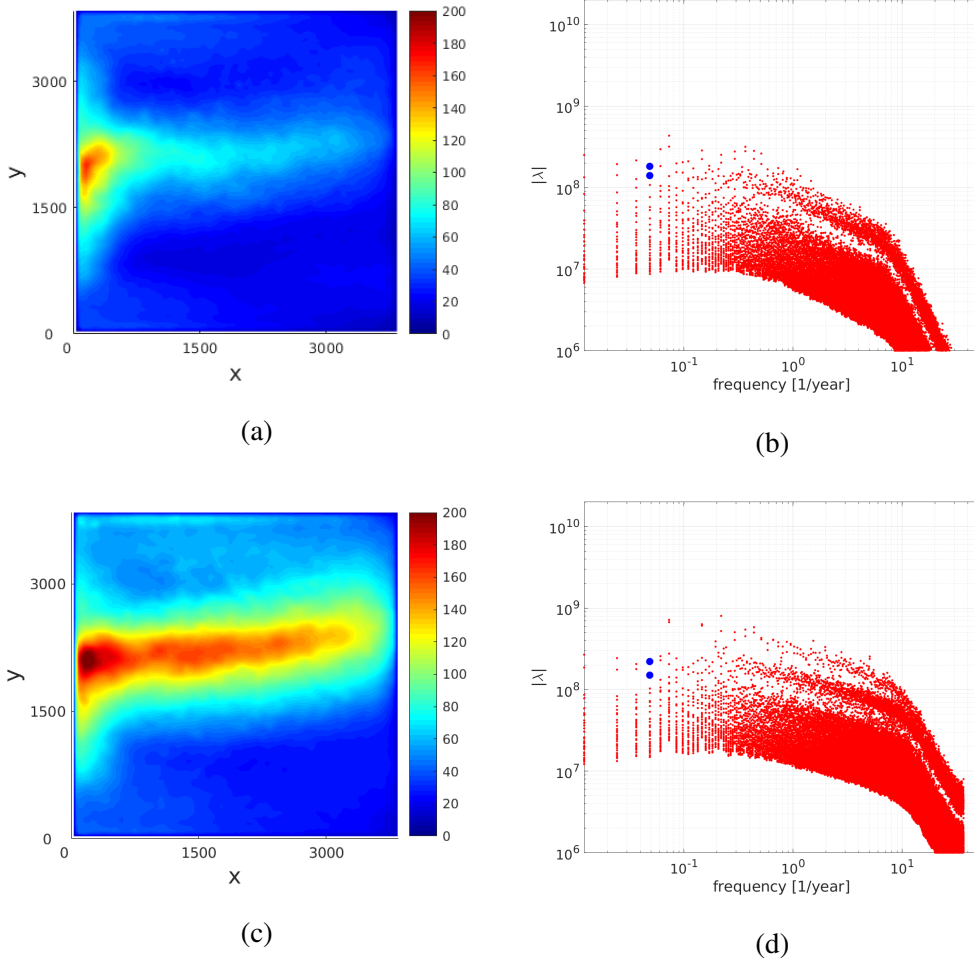


Figure 3: Statistics of the upper-layer PV anomaly field for the low-resolution augmented solution (129^2 grid) obtained by feeding the true eddy field extracted with the $W = 21$ filter): (a) standard deviation showing partial reconstruction of the eastward jet extension; (b) temporal spectral content provided by DAHD; the LFV (blue dots) is not reproduced, compared to the reference truth in Fig. 2b. Panels (c)-(d) are same as (a)-(b), but for the eddies extracted with the filter size $W = 41$; the eastward jet extension is now well reproduced, but there is still no LFV.

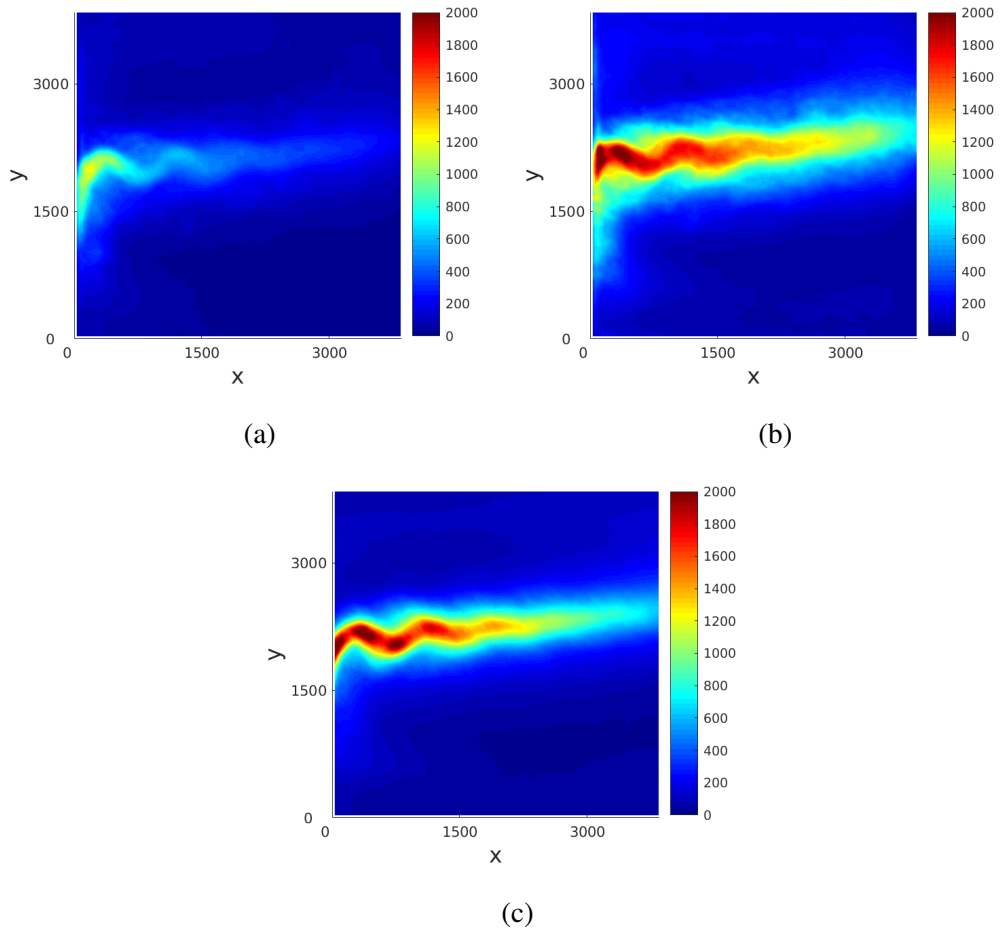


Figure 4: Standard deviations of different eddy forcings: (a) on-line eddy forcing from the solution augmented with eddies extracted with filter size $W = 21$; (b) same as (a), but for $W = 41$; (c) true (offline) eddy forcing, as in Ryzhov et al. (2019)). The on-line eddy forcing in (a) is about 4 times weaker than the off-line forcing, which is one of the reasons for the augmentation failure.

298 the filtering. The corresponding set of EOFs are obtained through the standard
 299 singular value decomposition, such that

$$\overline{Q}_{HR}^i = PC^i \cdot EOF^i, \quad (11)$$

300 where \overline{Q}_{HR}^i is the large-scale true PV anomaly in the i -th layer and in the matrix
 301 form rearranged so, that the rows correspond to the spatial degrees of freedom,
 302 whilst the columns represent their time evolutions; $PC^i = U^i \cdot S^i$, $EOF^i =$
 303 $(V^i)^*$, where U^i , S^i , V^i are the left eigenvector, diagonal singular value, and the
 304 right eigenvector matrices, respectively; \cdot^* is matrix transpose.

305 Projection of the on-line augmented PV anomaly Q^i onto some n EOFs EOF_n^i
 306 takes the form:

$$\overline{Q}_n^i = Q^i \cdot (EOF_n^i)^* \cdot EOF_n^i, \quad (12)$$

307 and the updated field \overline{Q}^i is used on the next time step of the model (Eq. 10).

308 There are two key parameters at the projection step: the number n of EOFs
 309 and the time interval T_{proj} between successive projections; these parameters are
 310 chosen empirically, for optimizing both the results and computational costs. We
 311 found by sensitivity experiments that the number of the EOFs should be relatively
 312 large, and 2000 out of $129^2 = 16641$ total EOFs are good enough; and T_{proj} should
 313 not be much longer than 100 model days, used here as the benchmark value. With
 314 these parameters, the augmented model recovered not only more than 95% of the
 315 LFV spectral power but also the correct frequencies. We varied the number of
 316 the EOFs and obtained qualitatively similar results within the 500–2000 range,
 317 and the lower values degrade the solution. Since the EOF projections are made
 318 infrequently, the filtering process is computationally inexpensive.

319 The additionally filtered model solutions now exhibit the LFV as diagnosed
 320 by DAHD spectra shown in Fig. 5a for $W = 21$ and Fig. 5b for $W = 41$. It is

321 worth noting that even in the solution augmented with weaker eddies ($W = 21$)
 322 the LFV is also reproduced, albeit it is not as energetic as with the stronger eddies
 323 ($W = 41$). The eastward jet extension is also reproduced similarly to the case
 324 without large-scale filtering (see Fig. 3).

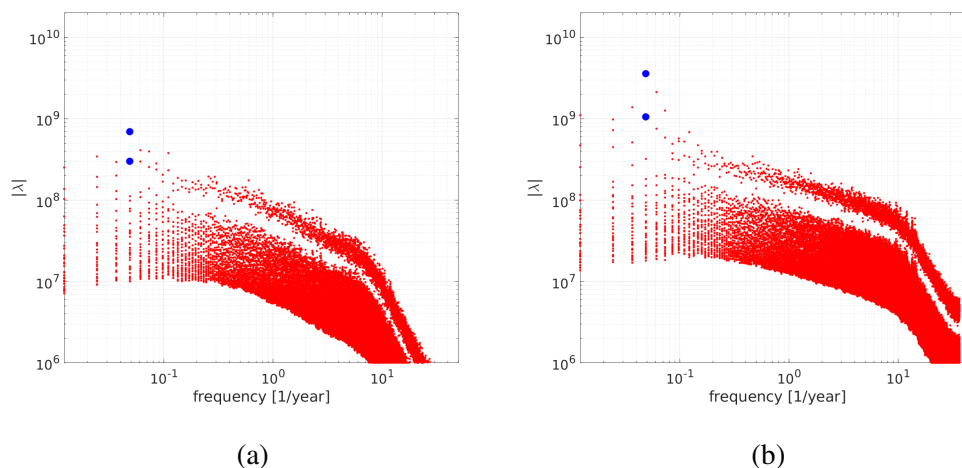


Figure 5: The DAHD temporal spectra of the upper-layer PV anomaly field in augmented and additionally filtered model solutions: (a) $W = 21$ (weaker eddies) ; (b) $W = 41$ (stronger eddies). The LFV (see the peaks with the blue dots) is now present in both solutions, and it is more intensive with stronger eddies.

325 In addition to the detailed DAHD spectral space-time diagnostic of PV anomaly
 326 field, it is also useful to consider the manifestation of LFV in the total poten-
 327 tial energy, which is a global characteristic of the solution. Figure 6 shows
 328 the Fourier spectral analysis of the potential energy time series by the standard
 329 Multitaper method (Percival and Walden, 1993), which reveals broadband LFV
 330 peaks at frequency $\approx 0.06 \text{ year}^{-1}$ (about 17 years period), both for the refer-
 331 ence high-resolution and augmented low-resolution solutions, whilst the reference
 332 low-resolution solution features no LFV with a mostly flat spectrum. Due to the

333 projection, the augmented solution acquires oversaturated high frequencies near
 334 the LFV peak; this may be dealt with by carefully selecting the projection basis of
 335 the filtering procedure so to filter out spurious small-scale effects and is beyond the
 336 scope of the current study as we aimed at imbuing the coarse-resolution solution
 337 with the correct LFV.

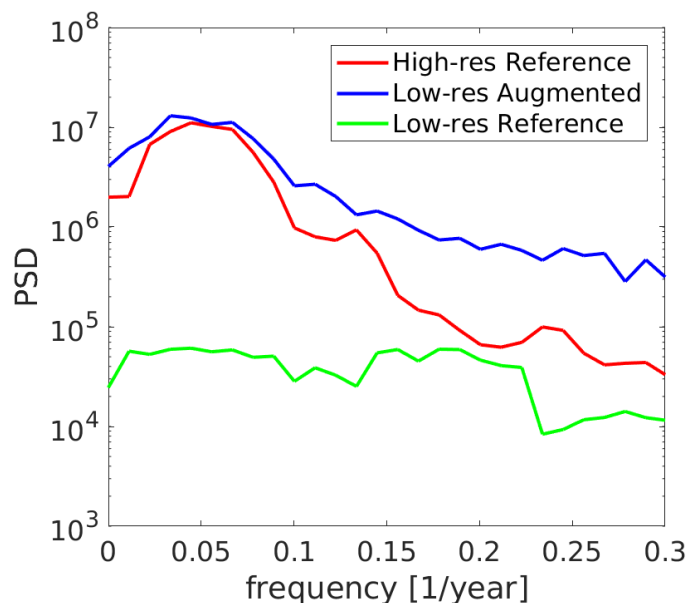


Figure 6: Power spectrum density (PSD) of the potential energy by the Multitaper method, featuring the energetic and broadband LFV with the main period of ≈ 17 years, in both the reference high-resolution solution and augmented low-resolution solution (supplied by the eddy field obtained with filter $W = 41$ and periodically projected onto 2000 EOFs of the large-scale "truth" basis), as opposed to the lack of such LFV in the reference low-resolution solution.

338 Finally, we would like to emphasize that feeding the eddies to induce the aug-
 339 menting eddy forcing in the low-resolution model (Eq. 9) is absolutely necessary
 340 for generating the LFV, and we verified this by turning it off. If the filtering based
 341 on the EOF projection procedure is applied alone, it does not augment the solu-

342 tion thus confirming that the main component of the parameterization is the eddy
343 forcing.

344 **5. Statistical emulation of the eddy field**

345 Here we developed data-driven statistical emulators of the true eddy field for
346 feeding them into the low-resolution model instead of the original high-resolution
347 eddy fields. The number of statistical emulation methods has recently surged,
348 including stochastic approaches in climate science (Penland and Matrosova, 2001;
349 Strounine et al., 2010; Franzke et al., 2015; Kondrashov et al., 2015; Chen et al.,
350 2016; Palmer, 2019; Seleznev et al., 2019; Foster et al., 2020), as well as other
351 machine-learning (deep learning) methods developed for fluid dynamics appli-
352 cations (Brunton et al., 2020; Bolton and Zanna, 2019). The detailed analysis
353 of emulated eddy fields is beyond the scope of this study, and in the context of
354 assessing the skill of our emulators we focus solely on one of the central problems
355 in climate ocean model simulations, namely, the correct rectification of the eddy
356 field’s impact on the large-scale circulation. Thus we aimed for the solution of
357 the low-resolution model, when augmented by an emulated eddy field, to be able
358 to reproduce the long-term statistics of the high-resolution reference solution. We
359 utilized the same skill measures as for the true eddy field explored in previous
360 section. These are the geometrical shape of the large-scale circulation patterns, as
361 well as the manifestation of the LFV.

362 We used a 30000-day long high-resolution dataset of the eddy stream function
363 $\hat{\Psi}$ for the three layers combined. The dataset is then coarse-grained onto the low
364 spatial resolution (129×129), and further compressed by the PCA. We retained
365 the leading 1000 PCs that account for $\approx 98\%$ of the variability.

366 As a basic and most straightforward emulator, we considered a linear stochastic
 367 regression model (Kravtsov et al., 2005, 2006; Kondrashov et al., 2005, 2015) in
 368 the following discrete form:

$$\xi_{t+1} - \xi_t = \mathbf{A}\xi_t + r_t^{(0)}, \quad (13)$$

369 where t is the time index (in days), ξ is a vector of PCs, and \mathbf{A} is a matrix of
 370 the regression coefficients. While Eq. 13 can include additional model layers of
 371 hidden variables obtained in a sequential regression procedure, it is not necessary
 372 here since the regression residual $r_t^{(0)}$ is well approximated by a spatially correlated
 373 white noise, $r_t^{(0)} = \Sigma \dot{W}$, where W is a Wiener process and Σ is the Cholesky
 374 decomposition of the correlation matrix of the residuals from the model fitting.

375 The emulated PCs are obtained by initializing the model from the first data
 376 point of the training interval and by running it for 30000 days. The eddy field is
 377 reconstructed in space from the emulated PCs by using the EOF basis, and then
 378 it is fed into the low resolution model in our augmentation procedure. While this
 379 basic emulator of the eddy field yields a fairly reasonable geometrical structure of
 380 the jet extension in the augmented solution (Fig. 7a), it does not induce the LFV
 381 as evident by the flat spectral density curve of the full potential energy (Fig. 7b),
 382 which is also similar to the non-augmented low-resolution solution.

383 A closer analysis shows that the lack of the LFV in the augmented solution
 384 is related to the spectral content of the emulated eddy field, in which energy at
 385 low frequencies is underestimated in comparison to the true eddy field. In turn,
 386 because the LFV in the true eddy field is considerably weaker than in the true
 387 reference solution, it is challenging to capture it by an emulator based on PCA
 388 PCs, which typically mix different temporal scales.

389 The DAHD method (sec. 2.2 and Appendix A) provides a novel emulation

390 alternative, as it combines identification of frequency-ranked modes and their
 391 efficient modelling. It extracts pairs of data-adaptive harmonic modes (DAHMs)
 392 that form an orthonormal set of spatial patterns oscillating harmonically in time,
 393 and, thus, represent global monochromatic space-time filters. Projection of the
 394 dataset onto DAHMs yields pairs of narrowband time series of data-adaptive
 395 harmonic coefficients (DAHCs), which are modulated in amplitude, but do not
 396 mix temporal scales.

397 Chekroun and Kondrashov (2017) showed that the Stuart–Landau (SL) stochas-
 398 tic oscillator – a nonlinear oscillating system near a Hopf bifurcation and driven
 399 by an additive noise, is best suited to model amplitude modulations and fre-
 400 quency for the narrowband and in-phase quadrature time series of a DAHC pair
 401 $(\zeta_t^+(f), \zeta_t^-(f))$, associated with a given spectral pair $(\lambda^+(f), \lambda^-(f))$ (see Sec. 2.2
 402 and Appendix B), here written in a compact form with a complex number notation:

$$z_{t+1}(f) - z_t(f) = (\mu(f) + i\gamma(f))z_t(f) - (1 + i\beta(f))|z_t(f)|^2 z_t(f) + \epsilon_t, \quad (14)$$

403 where $z_t(f) = \zeta_t^+(f) + i\zeta_t^-(f)$, $\mu(f)$, $\gamma(f)$ and $\beta(f)$ are real parameters and ϵ_t is
 404 an additive noise. Furthermore, multiple SL-oscillators associated with the same
 405 non-zero frequency are linearly coupled and synchronized across frequencies by
 406 the pairwise-correlated white noise, while the model parameters are estimated by
 407 a regression with constraints (see Appendix B for numerical details). The original
 408 dataset with its multiple time scales can be modeled in a computationally efficient
 409 manner since the contribution of each temporal frequency is simulated in parallel.

410 Kondrashov et al. (2018) developed a stochastic DAHD emulator for the LFV
 411 in the model considered, and here we extended these results to the eddies. We
 412 used the leading $d = 100$ PCs of the eddy streamfunction capturing $\approx 70\%$ of the
 413 variance and applied the DAHD with the embedding window of $M = 100$ days.

414 Then, we fit the model of coupled $d = 100$ stochastic oscillators for the DAHCs
 415 and obtained their emulations for the $M = 100$ frequencies. After emulated
 416 DAHCs were back-transformed into the space-time eddy field by using DAHMs
 417 and EOFs, and combined across all the emulated frequencies, we fed the outcome
 418 into the augmented model. The geometrical shape of the augmented solution is
 419 again reproduced fairly well, and it is very similar to Fig. 7a (not shown for
 420 brevity). Furthermore, since the LFV is now better captured in the emulated eddy
 421 field (compare to the high-resolution "truth"), it is also induced in the augmented
 422 solution (Fig. 7b), albeit it is less energetic then when the true eddy field is used
 423 (see Fig. 6).

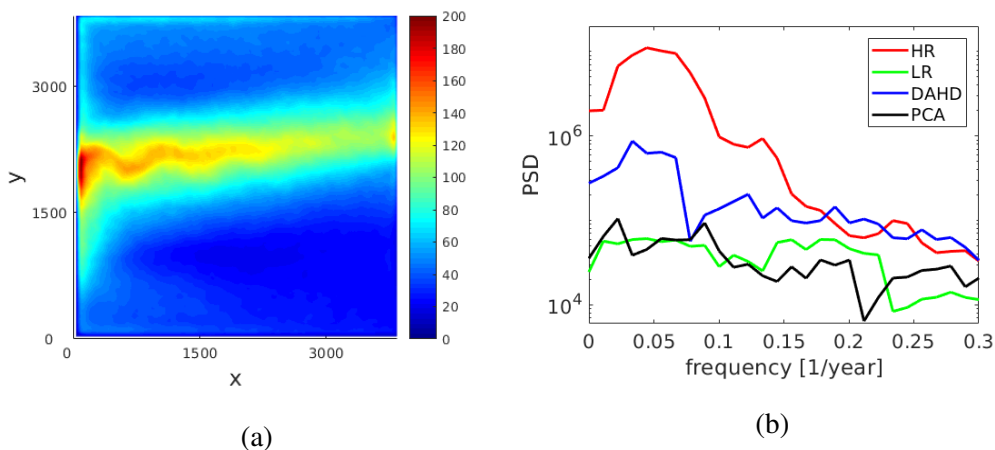


Figure 7: (a) Standard deviation of the upper-layer PV anomalies in the augmented solution with an artificial eddy field emulated by a PCA-based linear model (Eq. 13) (the periodical projection onto the 2000 EOFs of the large-scale "truth" basis is applied as well). The pattern of the standard deviation for the case of the DAHD model (Eq. 14) is similar (however, its magnitude is noticeably larger) and is not shown for brevity; (b) Power spectrum densities of the potential energy by the Multitapering method: the LFV is reproduced much better in the case of the DAHD-emulated eddy field.

424 **6. Conclusions**

425 In this paper we focused on improving solutions of an eddy-permitting low-
426 resolution model by augmenting it with the information from the reference high-
427 resolution model solution, which was treated as the observed truth. Our approach
428 can be viewed as a basis for developing data-driven parameterizations for the
429 mesoscale oceanic eddies and their effects, and in perspective for other types of
430 turbulent fluid motions. Ultimately, the parameterization should involve statistical
431 emulations of the key unresolved or under-resolved flow features. We adopted a
432 systematic approach towards such a parameterization framework; this paper is the
433 second one in the series, after (Ryzhov et al., 2019).

434 For the ocean circulation model, we considered the classical, wind-driven
435 double gyres in the quasigeostrophic approximation with 3 active isopycnal layers,
436 and in an idealized, closed, midlatitude basin configuration. Solutions of the
437 double-gyre model are notoriously sensitive to the spatial grid resolution, which
438 is typical for the general ocean circulation models. Two prominent flow features,
439 which are crucially dependent on the resolution, are in the focus of our study:
440 (1) the eastward jet extension of the western boundary currents with its adjacent
441 recirculation zones, and (2) the intrinsic, large-scale low-frequency (interdecadal)
442 variability of the gyres that is most pronounced in the eastward jet region. Both of
443 these features are essentially mesoscale eddy-driven, therefore, for their dynamical
444 representation in the model the eddies have to be either properly resolved, which
445 is computationally expensive, or adequately parameterized in terms of a simpler
446 model.

447 In the high-resolution reference solution both of the key features are well
448 represented, whereas the low-resolution reference solution lacks any of them. Mo-

449 tivation for including (1) is straightforward, because any eddy parameterization
450 is, first of all, tested for its ability to simulate the large-scale climatological fields.
451 Motivation for including (2) is to test the ability of the parameterization to simulate
452 intrinsic climate variabilities similar to the relatively well understood interdecadal
453 variability featured in our model. Our hope is that testing mesoscale eddy param-
454 eterization skills will eventually include climate variability signals as the standard
455 test beds.

456 Our model augmentation procedure involves the following main steps. First,
457 the high-resolution (true) solution is decomposed into large-scale and small-scale
458 (eddy) flow components by simple moving-average filtering in space. This flow
459 decomposition is neither unique nor obviously constrained by dynamical or sta-
460 tistical arguments. Here, we only assumed that the filter width should be about
461 scaled with the first baroclinic Rossby deformation radius, since our study targets
462 mesoscale eddies.

463 In the prequel study (Ryzhov et al., 2019), the decomposed flow components
464 were used to find the history of the eddy forcing, which is just part of the advection
465 operator that involves the eddy field; then, this history was coarse-grained and
466 applied to augment the low-resolution model with many analyses and sensitivity
467 studies attached to this statement and reported in the paper. In the present study
468 we extended the approach by supplying the primary eddy fields instead of the
469 eddy forcing, which is a higher-level and subtler information. Moreover, we tested
470 the augmentation procedure skills in terms of the challenging reproduction of the
471 LFV. The eddy field component was interactively coupled with the corresponding
472 low-resolution model solution, which was treated as the simulated large-scale flow
473 component, via the (on-line) eddy forcing operator, which can be viewed as an

474 additional dynamical constraint imposed on the augmentation procedure.

475 We found that the augmentation significantly improved representation of the
476 eastward jet extension, but the LFV was still missing. The immediate hypothesis
477 was that this was because the eddies are too weak, hence, the interactive eddy
478 forcing was too weak to generate the LFV. We tested this hypothesis by increasing
479 the filter size used to extract the eddies, and the resulting new eddy forcing turned
480 out to be of the same intensity as the true eddy forcing; however, this further
481 improved the modelled eastward jet but did not generate the LFV. From this
482 we concluded that the LFV was crucially dependent on the correlations between
483 the large-scale flow and the eddy forcing, which were not fully respected by the
484 augmentation procedure.

485 We also realized that the eddy history alone was not sufficient, and some
486 additional information had to be supplied as part of the augmentation. We do not
487 yet have the ultimate answer on what this information should be, but in order to
488 make progress we decided to supply some large-scale flow information in terms of
489 interactive, weak filtering of the simulated large-scale flow towards the observed
490 truth. This idea was implemented as a statistical filtration - interactively projecting
491 the simulated transient flow anomalies onto the leading empirical orthogonal
492 functions (EOFs) of the reference (high-resolution) true flow.

493 This approach worked well, and we experimentally found the optimal number
494 of the EOFs and the optimal frequency of the applied filtering procedure, so that
495 the LFV was almost fully recovered. Since the filtering can be applied infrequently
496 (about every 100 days in our case) rather than continuously, which is also pos-
497 sible, its computational cost is nearly negligible. However, the exact amount of
498 information needed from the high-resolution "truth" for a correct rectification of

499 the LFV remains unknown and its assessment should be addressed elsewhere. We
500 hypothesised that this information should contain correct correlations between the
501 eddy and large-scale fields. We also demonstrated that the filtering was of sec-
502 ondary importance relative to the supplied eddy forcing, because when the latter
503 was switched off, the filtering alone was not capable of augmenting the solution
504 to any acceptable level.

505 Finally, we developed a statistical emulation of the eddy field as spatio-temporal
506 stochastic process, and used it in our augmented procedure. Results showed that
507 the frequency-ranked data-adaptive harmonic decomposition (DAHD) emulator re-
508 produces the LFV substantially better than the PCA-based linear stochastic model.

509 An agenda for further research stemming from this paper is to build on and
510 improve statistical emulators for the eddy field, as well as to consider extending the
511 proposed approach beyond the relatively simple quasigeostrophic approximation
512 to comprehensive general circulation models. Constraining the large-scale/eddy
513 flow decomposition and making it consistent with the low-resolution ocean model
514 is also very important. Finally, adding new criteria (e.g., higher-order statistical
515 moments and spatio-temporal correlations) for assessing eddy parameterization
516 skills should not be too far away.

517 **Acknowledgements**

518 We thank Dr J. Maddison and one anonymous reviewer for constructive com-
519 ments that helped improve this manuscript. This research was supported by the
520 National Science Foundation (NSF) grants *OCE* – 1658357 and the NERC grant
521 *NE/R011567/1*. Pavel Berloff also gratefully acknowledges funding by NERC
522 Grant No. *NE/T002220/1* and Leverhulme Grant No. *RPG* – 2019 – 024.

523 DAHD analysis was supported by the Russian Science Foundation (Grant No.
524 18 – 12 – 00231). We would like to acknowledge the high-performance computing
525 support from Cheyenne (doi:10.5065/D6RX99HX) provided by NCAR’s Com-
526 putational and Information Systems Laboratory, sponsored by the NSF. The DAHD
527 Toolbox is available at: <http://research.atmos.ucla.edu/tcd/dkondras/Software.html>

528 **Appendix A. Data-adaptive harmonic decomposition (DAHD)**

529 Here we present a brief summary of the DAHD frequency-domain imple-
530 mentation and stochastic emulation methodology following (Chekroun and Kon-
531 drashov, 2017; Kondrashov and Chekroun, 2018; Kondrashov et al., 2018,b) and
532 tailored to high-dimensional datasets. We consider a multivariate time series
533 $\mathbf{X}(t) = (X_1(t), \dots, X_d(t))$ formed with d spatial channels and $t = 1, \dots, N$ time
534 points (sampled evenly). Double-sided (unbiased) cross-correlation coefficients
535 $\rho^{(p,q)}(m)$ are estimated for all the pairs of channels p and q and time lag m up to a
536 maximum $M - 1$:

$$\rho^{(p,q)}(m) = \begin{cases} \frac{1}{N-m} \sum_{t=1}^{N-m} X_p(t+m)X_q(t), & 0 \leq m \leq M - 1, \\ \rho^{(q,p)}(-m), & m < 0. \end{cases} \quad (15)$$

537 where M is the embedding window and each of $\rho^{(p,q)}(m)$ sequences is of length
538 $M' = 2M - 1$. The DAHD numerical algorithm computes its spectral elements
539 $(\lambda_j, \mathbf{W}_j, j = 1, \dots, d(2M - 1))$ by utilizing a $d \times d$ symmetrized complex cross-
540 spectral matrix $\mathfrak{S}(f)$ built from the Fourier transforms of the cross-correlation
541 sequences (see Eq. 4). The data-adaptive harmonic modes (DAHMs) represent
542 collection of spatio-temporal patterns $\mathbf{W}_j = (\mathbf{E}_1^j, \dots, \mathbf{E}_d^j)$ oscillating with differ-

543 ent but single frequency f in time-embedded space $1 \leq m \leq M'$:

$$\mathbf{E}_k^j(m) = B_k^j \cos(2\pi f m + \theta_k^j), \quad 1 \leq k \leq d, \quad (16)$$

544 where the amplitudes B_k^j and phases θ_k^j are data-adaptive, f takes distinct M' values
545 that are equally spaced in Nyquist interval $[0 \ 0.5]$,

$$f = \frac{(\ell - 1)}{M' - 1}, \quad \ell = 1, \dots, \frac{M' + 1}{2}, \quad (17)$$

546 and $|\lambda_j|$ informs on energy conveyed by \mathbf{W}_j . In particular, for each $f \neq 0$,
547 there are $2d$ positive-negative eigenelements which are necessarily paired as
548 $(\lambda_k^+(f) = -\lambda_k^-(f), k = 1, \dots, d)$, while the phases for the associated DAHM
549 pair $(\mathbf{W}_k^+(f), \mathbf{W}_k^-(f))$ satisfy $\theta_k^+ = \theta_k^- + \pi/2$, i.e. these modes are shifted by one
550 fourth of the period and are thus always in exact phase quadrature, similar to the
551 sine-and-cosine pair in the Fourier analysis, but in a data-adaptive and global-in-
552 space fashion. There are also d (non paired) spectral elements $(\lambda_k, \mathbf{W}_k)$ associated
553 with the frequency $f = 0$. The Fourier transforms of the DAHMs are computed as
554 eigenvectors of the matrix $\mathfrak{S}(f)\overline{\mathfrak{S}(f)}$ (Chekroun and Kondrashov, 2017, Theorem
555 V.1 and Eq.74):

$$\mathfrak{S}(f)\overline{\mathfrak{S}(f)}\widehat{W}_k(f) = \lambda_k^2\widehat{W}_k(f) \quad (18)$$

556 and spatiotemporal patterns of $(\mathbf{W}_k^+(f), \mathbf{W}_k^-(f))$ are obtained then by the inverse
557 Fourier transform. A projection of X onto given \mathbf{W}_j yields the time series of the
558 DAHD expansion coefficients (DAHCs):

$$\zeta_j(t) = \sum_{m=1}^{M'} \sum_{k=1}^d X_k(t + m - 1) \mathbf{E}_k^j(m) \quad (19)$$

559 where $1 \leq t \leq N - M' + 1$. The time series of a given DAHC pair $(\zeta_k^+(t), \zeta_k^-(t))$
560 associated with the modes $(\mathbf{W}_k^+(f), \mathbf{W}_k^-(f))$ at the frequency $f \neq 0$, are narrow-
561 band, nearly in phase quadrature and heavily modulated in amplitude.

562 **Appendix B. Frequency-Ranked Stochastic Emulators**

563 The collective behavior of the d pairs at the frequency $f \neq 0$ (see Appendix A)
 564 is simulated by a system of linearly coupled Stuart-Landau stochastic oscillators:

$$\begin{aligned}
 \frac{d\zeta_k^+}{dt} &= \beta_k(f)\zeta_k^+ - \alpha_k(f)\zeta_k^- - \sigma_k(f)\zeta_k^+((\zeta_k^+)^2 + (\zeta_k^-)^2) \\
 &\quad + \sum_{i \neq k}^d a_{ik}(f)\zeta_i^+ + \sum_{i \neq k}^d b_{ik}(f)\zeta_i^- + \epsilon_k^+, \\
 \frac{d\zeta_k^-}{dt} &= \alpha_k(f)\zeta_k^+ + \beta_k(f)\zeta_k^- - \sigma_k(f)\zeta_k^-((\zeta_k^+)^2 + (\zeta_k^-)^2) \\
 &\quad + \sum_{i \neq k}^d c_{ik}(f)\zeta_i^+ + \sum_{i \neq k}^d d_{ik}(f)\zeta_i^- + \epsilon_k^-,
 \end{aligned} \tag{20}$$

565 where $1 \leq k \leq d$; the model parameters are estimated by a pairwise multiple linear
 566 regression with linear constraints on $\alpha_k(f)$ and $\beta_k(f)$ to ensure antisymmetry
 567 for the linear coupling within a given pair, as well as equal and positive values
 568 $\sigma_k(f) > 0$ to ensure numerical stability. The stochastic forcing in Eq. 20 is
 569 informed by regression residuals from the model fitting, namely $\begin{bmatrix} \epsilon_t^+ \\ \epsilon_t^- \end{bmatrix} = \Sigma d\mathbf{W}$,
 570 where Σ is the $2d \times 2d$ Cholesky decomposition of the correlation matrix of the
 571 residuals and $d\mathbf{W}$ is a $2d$ -valued Wiener process. The linear stochastic emulator
 572 (Eq. 13) is used to model the time series of the DAHCs associated with $f \equiv 0$,
 573 which are not paired.

574 Any subset of DAHCs can be convolved with its corresponding set of DAHMs,
 575 to produce a partial or full reconstruction of the original dataset. Thus, the
 576 following j -th reconstructed component (RC) at time t and for channel k is defined
 577 as:

$$R_k^j(t) = \frac{1}{M_t} \sum_{m=L_t}^{U_t} \zeta_j(t-m+1) \mathbf{E}_k^j(m), \quad 1 \leq m \leq M' \tag{21}$$

578 where $L_t(U_t)$ is a lower (upper) bound in $\{1, \dots, M'\}$ that depends on time and
579 the normalization factor M_t equals M' except near the ends of the time series.
580 The sum of all the RCs across all the frequencies recovers the original time series,
581 and stochastically emulated DAHCs are back-transformed to the phase-space of
582 the original dataset by using Eq. 21.

583 **References**

- 584 I. V. Shevchenko, P. S. Berloff, Multi-layer quasi-geostrophic ocean dynamics in
585 eddy-resolving regimes, *Ocean Model.* 94 (2015) 1–14.
- 586 I. Shevchenko, P. Berloff, D. Guerrero-Lopez, J. Roman, On low-frequency vari-
587 ability of the midlatitude ocean gyres, *J. Fluid Mech.* 795 (2016) 423–442.
- 588 D. P. Marshall, J. R. Maddison, P. S. Berloff, A framework for parameterizing
589 eddy potential vorticity fluxes, *J. Phys. Oceanogr.* 42 (2012) 539–557.
- 590 S. D. Bachman, B. Fox-Kemper, B. Pearson, A scale-aware subgrid model for
591 quasi-geostrophic turbulence, *J. Geophys. Res.: Oceans* 122 (2017) 1529–1554.
- 592 P. R. Gent, J. C. McWilliams, Isopycnal mixing in ocean circulation models, *J.*
593 *Phys. Oceanogr.* 20 (1990) 150–155.
- 594 J. S. Frederiksen, Subgrid-scale parameterizations of eddy-topographic force,
595 eddy viscosity, and stochastic backscatter for flow over topography, *J. Atmos.*
596 *Sci.* 56 (1999) 1481–1494.
- 597 J. S. Frederiksen, T. J. O’Kane, M. J. Zidikheri, Stochastic subgrid parameteriza-
598 tions for atmospheric and oceanic flows, *Physica Scripta* 85 (2012) 068202.

- 599 P. Porta Mana, L. Zanna, Toward a stochastic parametrization of ocean mesoscale
600 eddies, *Ocean Model.* 79 (2014) 1–20.
- 601 P. Berloff, Dynamically consistent parameterization of mesoscale eddies. Part I:
602 Simple model, *Ocean Modelling* 87 (2015) 1 – 19.
- 603 P. Berloff, Dynamically consistent parameterization of mesoscale eddies. Part II:
604 Eddy fluxes and diffusivity from transient impulses, *Fluids* 1 (2016).
- 605 L. Zanna, P. Porta Mana, J. Anstey, T. David, T. Bolton, Scale-aware deterministic
606 and stochastic parametrizations of eddy-mean flow interaction, *Ocean Model.*
607 111 (2017) 66–80.
- 608 P. Berloff, Dynamically consistent parameterization of mesoscale eddies. Part III:
609 Deterministic approach, *Ocean Modelling* 127 (2018) 1 – 15.
- 610 J. Mak, J. R. Maddison, D. P. Marshall, D. R. Munday, Implementation of a
611 geometrically informed and energetically constrained mesoscale eddy parame-
612 terization in an ocean circulation model, *Journal of Physical Oceanography* 48
613 (2018) 2363–2382.
- 614 E. Ryzhov, D. Kondrashov, N. Agarwal, P. Berloff, On data-driven augmentation
615 of low-resolution ocean model dynamics, *Ocean Modelling* 142 (2019) 101464.
- 616 A. J. Majda, I. Timofeyev, E. Vanden Eijnden, Models for stochastic climate
617 prediction, *PNAS* 96 (1999) 14687–14691.
- 618 I. Fatkullin, E. Vanden-Eijnden, A computational strategy for multiscale systems
619 with applications to lorenz 96 model, *Journal of Computational Physics* 200
620 (2004) 605 – 638.

- 621 S. Kravtsov, D. Kondrashov, M. Ghil, Multi-level regression modeling of nonlinear
622 processes: Derivation and applications to climatic variability, *J. Climate* 18
623 (2005) 4404–4424.
- 624 D. Crommelin, E. Vanden-Eijnden, Subgrid-scale parameterization with condi-
625 tional markov chains, *Journal of the Atmospheric Sciences* 65 (2008) 2661–
626 2675.
- 627 H. M. Arnold, I. M. Moroz, T. N. Palmer, Stochastic parametrizations and
628 model uncertainty in the Lorenz 96 system, *Philosophical Transactions of
629 the Royal Society A: Mathematical, Physical and Engineering Sciences* 371
630 (2013) 20110479.
- 631 A. J. Chorin, F. Lu, Discrete approach to stochastic parametrization and dimension
632 reduction in nonlinear dynamics, *PNAS* 112 (2015) 9804–9809.
- 633 J. R. Maddison, D. P. Marshall, J. Shipton, On the dynamical influence of ocean
634 eddy potential vorticity fluxes, *Ocean Modelling* 92 (2015) 169 – 182.
- 635 Y. K. Ying, J. R. Maddison, J. Vanneste, Bayesian inference of ocean diffusivity
636 from lagrangian trajectory data, *Ocean Modelling* 140 (2019) 101401.
- 637 S. Karabasov, P. Berloff, V. Goloviznin, CABARET in the ocean gyres, *Ocean
638 Modelling* 30 (2009) 155 – 168.
- 639 M. F. Jansen, I. M. Held, Parameterizing subgrid-scale eddy effects using ener-
640 getically consistent backscatter, *Ocean Modelling* 80 (2014) 36 – 48.
- 641 M. F. Jansen, I. M. Held, A. Adcroft, R. Hallberg, Energy budget-based backscatter

- 642 in an eddy permitting primitive equation model, *Ocean Modelling* 94 (2015)
643 15 – 26.
- 644 I. Shevchenko, P. Berloff, Eddy backscatter and counter-rotating gyre anomalies
645 of midlatitude ocean dynamics, *Fluids* 1 (2016).
- 646 P. S. Berloff, J. McWilliams, Large-scale, low-frequency variability in wind-driven
647 ocean gyres, *J. Phys. Oceanogr.* 29 (1999) 1925–1949.
- 648 M. D. Chekroun, D. Kondrashov, Data-adaptive harmonic spectra and multilayer
649 Stuart-Landau models, *Chaos* 27 (2017) 093110.
- 650 D. Kondrashov, M. D. Chekroun, P. Berloff, Multiscale Stuart-Landau emulators:
651 Application to wind-driven ocean gyres, *Fluids* 3 (2018) 21.
- 652 D. Kondrashov, E. A. Ryzhov, P. Berloff, Data-adaptive harmonic analysis of
653 oceanic waves and turbulent flows, *Chaos: An Interdisciplinary Journal of*
654 *Nonlinear Science* 30 (2020) 061105.
- 655 D. Kondrashov, M. D. Chekroun, X. Yuan, M. Ghil, Data-Adaptive Harmonic
656 Decomposition and Stochastic Modeling of Arctic Sea Ice, in: A. A. Tsonis
657 (Ed.), *Advances in Nonlinear Geosciences*, Springer International Publishing,
658 Cham, 2018a, pp. 179–205. doi:10.1007/978-3-319-58895-7_10.
- 659 D. Kondrashov, M. D. Chekroun, M. Ghil, Data-adaptive harmonic decomposition
660 and prediction of Arctic sea ice extent, *Dynamics and Statistics of the Climate*
661 *System* 3 (2018b).
- 662 D. Kondrashov, M. D. Chekroun, Data-adaptive harmonic analysis and modeling

- 663 of solar wind-magnetosphere coupling, *Journal of Atmospheric and Solar-*
664 *Terrestrial Physics* (2018).
- 665 R. W. Preisendorfer, *Principal Component Analysis in Meteorology and Oceanog-*
666 *raphy*, Elsevier, New York, 425 pp., 1988.
- 667 P. Berloff, A. Hogg, W. Dewar, The turbulent oscillator: A mechanism of low-
668 frequency variability of the wind-driven ocean gyres, *J. Phys. Oceanogr.* 37
669 (2007) 2363–2386.
- 670 D. Kondrashov, P. Berloff, Stochastic modeling of decadal variability in ocean
671 gyres, *Geophys. Res. Lett.* 42 (2015) 1543–1553.
- 672 K. Hasselmann, PIPs and POPs: The reduction of complex dynamical systems
673 using principal interaction and oscillation patterns, *Journal of Geophysical*
674 *Research: Atmospheres* 93 (1988) 11015–11021.
- 675 H. von Storch, G. Bürger, R. Schnur, J.-S. von Storch, Principal oscillation patterns:
676 A review, *Journal of Climate* 8 (1995) 377–400.
- 677 P. J. Schmid, Dynamic mode decomposition of numerical and experimental data,
678 *J. Fluid Mech.* 656 (2010) 5–28.
- 679 H. Li, J.-S. von Storch, On the fluctuating buoyancy fluxes simulated in a ogcm,
680 *J. Phys. Oceanogr.* 43 (2013) 1270–1287.
- 681 H. A. Dijkstra, *Nonlinear Climate Dynamics*, Cambridge University Press, Cam-
682 *bridge*, UK, 2013.
- 683 H. A. Dijkstra, A normal mode perspective of intrinsic ocean-climate variability,
684 *Annu. Rev. Fluid Mech.* 48 (2018) 341–363.

- 685 J. Viebahn, D. Crommelin, H. Dijkstra, Toward a turbulence closure based on
686 energy modes, *Journal of Physical Oceanography* 49 (2019) 1075–1097.
- 687 N. Agarwal, E. Ryzhov, D. Kondrashov, P. Berloff, Scale-aware flow decomposition
688 and statistical analysis of the eddy forcing, Submitted (2020).
- 689 P. Berloff, On dynamically consistent eddy fluxes, *Dyn. Atmos. Ocean.* 38 (2005)
690 123–146.
- 691 D. B. Percival, A. T. Walden, *Spectral analysis for physical applications*, Cambridge
692 university press, 1993.
- 693 C. Penland, L. Matrosova, Expected and Actual Errors of Linear Inverse Model
694 Forecasts, *Monthly Weather Review* 129 (2001) 1740–1745.
- 695 K. Strounine, S. Kravtsov, D. Kondrashov, M. Ghil, Reduced models of at-
696 mospheric low-frequency variability: Parameter estimation and comparative
697 performance, *Physica D: Nonlinear Phenomena* 239 (2010) 145 – 166.
- 698 C. L. E. Franzke, T. J. O’Kane, J. Berner, P. D. Williams, V. Lucarini, Stochastic
699 climate theory and modeling, *Wiley Interdiscip. Rev. Clim. Change* 6 (2015)
700 63–78.
- 701 D. Kondrashov, M. D. Chekroun, M. Ghil, Data-driven non-Markovian closure
702 models, *Physica D* 297 (2015) 33 – 55.
- 703 C. Chen, M. A. Cane, N. Henderson, D. E. Lee, D. Chapman, D. Kondrashov,
704 M. D. Chekroun, Diversity, Nonlinearity, Seasonality, and Memory Effect in
705 ENSO Simulation and Prediction Using Empirical Model Reduction, *Journal*
706 *of Climate* 29 (2016) 1809–1830.

- 707 T. N. Palmer, Stochastic weather and climate models, *Nature Reviews Physics* 1
708 (2019) 463–471.
- 709 A. Seleznev, D. Mukhin, A. Gavrilov, E. Loskutov, A. Feigin, Bayesian framework
710 for simulation of dynamical systems from multidimensional data using recurrent
711 neural network, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 29
712 (2019) 123115.
- 713 D. Foster, D. Comeau, N. M. Urban, A Bayesian Approach to Regional Decadal
714 Predictability: Sparse Parameter Estimation in High-Dimensional Linear In-
715 verse Models of High-Latitude Sea Surface Temperature Variability, *Journal of*
716 *Climate* 33 (2020) 6065–6081.
- 717 S. L. Brunton, B. R. Noack, P. Koumoutsakos, Machine learning for fluid mechan-
718 ics, *Annual Review of Fluid Mechanics* 52 (2020).
- 719 T. Bolton, L. Zanna, Applications of deep learning to ocean data inference and
720 subgrid parameterization, *Journal of Advances in Modeling Earth Systems* 11
721 (2019) 376–399.
- 722 S. Kravtsov, P. Berloff, W. Dewar, M. Ghil, J. McWilliams, Dynamical origin of
723 low-frequency variability in a highly nonlinear midlatitude coupled model, *J.*
724 *Climate* 19 (2006) 6391–6408.
- 725 D. Kondrashov, S. Kravtsov, A. W. Robertson, M. Ghil, A hierarchy of data-based
726 ENSO models, *Journal of Climate* 18 (2005) 4425–4444.