

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, P. Berloff^{a,e}

^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 of Russian Academy of Sciences, Vladivostok, 690041, Russia

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

^e Institute of Numerical Mathematics of the Russian Academy of Sciences, 119333, Moscow, Russia

ARTICLE INFO

Keywords:

Ocean dynamics
Mesoscale eddies
Eddy forcing
Parameterizations

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 field and the explicitly modeled 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 modeling 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.

[☆] All authors contributed equally to the manuscript preparation.

* Corresponding author at: Department of Mathematics, Imperial College London, London, SW7 2AZ, UK.

E-mail addresses: e.ryzhov@imperial.ac.uk (E.A. Ryzhov), 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 (P. Berloff).

<https://doi.org/10.1016/j.ocemod.2020.101664>

Received 5 March 2020; Received in revised form 24 June 2020; Accepted 8 July 2020

Available online 13 July 2020

1463-5003/© 2020 Elsevier Ltd. All rights reserved.

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

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

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

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

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

2. Double-gyre model

2.1. Governing equations

We use the same model configuration as in Ryzhov et al. (2019). The model has been extensively tested both in eddy-permitting and eddy-resolving regimes (Marshall et al., 2012; Maddison et al., 2015; Shevchenko and Berloff, 2015; Shevchenko et al., 2016; Ying et al., 2019). A brief description is as follows. The quasi-geostrophic (QG) potential vorticity (PV) evolution in 3 stacked isopycnal layers ($i = 1, 2, 3$ from top to bottom) with densities ρ_i ($\rho_1 = 1000$, $\rho_2 = 1001.498$, $\rho_3 = 1001.62$ kg m⁻³) and heights H_i ($H_1 = 250$, $H_2 = 750$, $H_3 = 3000$ m) is given by

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

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

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

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

where the asymmetry, tilt, and wind stress magnitude parameters are $A = 0.9$, $B = 0.2$, and $\tau_0 = 0.08$ N m⁻², respectively.

The PV anomalies and streamfunctions are related through

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

where the stratification parameters S_1 , S_{21} , S_{22} , S_3 are chosen to yield the first and second baroclinic Rossby deformation radii of 40 and 23 km, respectively. The boundary conditions are no-flow-through

and partial-slip (with the partial-slip length scale equal to 120 km); the mass is conserved in each layer. The model is solved using the high-resolution CABARET method that features a second-order, non-dissipative and low-dispersive, conservative advection scheme (Karabasov et al., 2009).

Given an adequately fine spatial resolution, the model is capable of resolving the eddies that maintain the well-developed eastward jet extension of the western boundary current. Otherwise, the eastward jet extension is under-predicted or even absent because the backscatter process of the energy transfer from the eddies to the large-scale flow is under-resolved by the model (Jansen and Held, 2014; Jansen et al., 2015; Shevchenko and Berloff, 2016; Berloff, 2018).

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

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

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

Not only the spatial patterns but also the temporal variabilities of the reference solutions are different. To reveal details of the latter, we used the Data-Adaptive Harmonic Decomposition (DAHD) method (Chekroun and Kondrashov, 2017; Kondrashov et al., 2018a), which characterizes a complex and multiscale spatio-temporal variability by extracting spatial data-adaptive harmonic modes (DAHMs) such that each one of them oscillates at a single temporal frequency and is spatially orthogonal to all other modes at that frequency (see Appendix A for details). The DAHD has been successfully applied to characterize variabilities in different geophysical datasets including ocean circulation (Kondrashov et al., 2018a; Ryzhov et al., 2019; Kondrashov et al., 2020), sea ice (Kondrashov et al., 2018c,b), and space physics (Kondrashov and Chekroun, 2018).

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

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

where $1 \leq p, q \leq d$; and $\widehat{\rho^{p,q}}(f)$ is the Fourier transform of the double-sided cross-correlation coefficients $\rho^{(p,q)}(m)$ estimated for all pairs of the channels (PCs) p and q , and for the time lag m , up to its maximum

$M - 1$; i.e. $-(M - 1) \leq m \leq M - 1$. Each singular value $\sigma_k(f)$ of $\mathfrak{S}(f)$ is associated with a pair of negative/positive eigenvalues ($\lambda_k^+(f)$, $\lambda_k^-(f)$) obtained by using the standard DAHD time-domain formulation and an eigen-decomposition of a matrix formed of the elements $\rho^{(p,q)}(m)$ (Kondrashov et al., 2018a; Ryzhov et al., 2019; Kondrashov et al., 2020):

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

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

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

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

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

The striking difference is the pronounced LFV in the high-resolution solution (see the blue dots in Fig. 2b at the period ≈ 17 years), and its complete absence in the low-resolution solution (Fig. 2a). This interdecadal LFV was studied elsewhere (Berloff and McWilliams, 1999; Berloff et al., 2007; Shevchenko et al., 2016), and here we just note that the quality of an augmented low-resolution model can be tested by the model's capability to simulate this LFV.

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

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

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

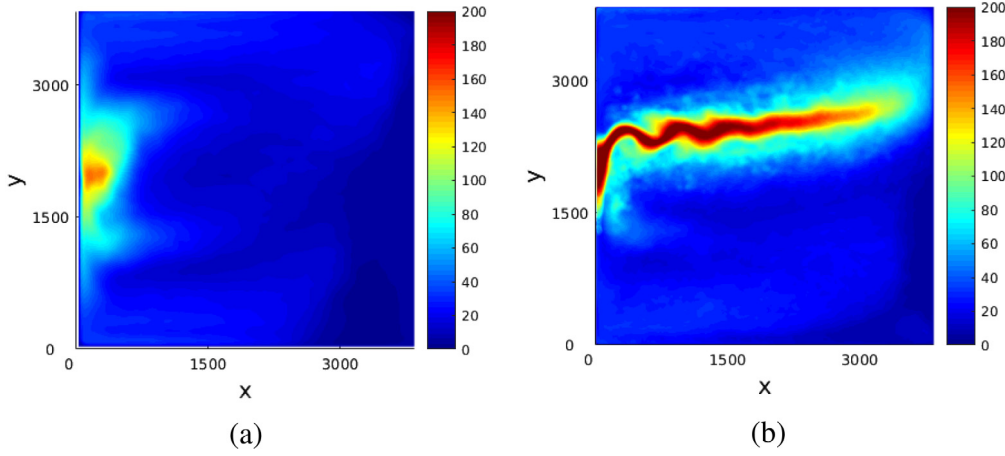


Fig. 1. Standard deviation of the upper-layer PV anomaly (q_i) produced by the (a) low-resolution (129^2) and (b) high-resolution (513^2) models. The solutions emphasize the crucial effect of the spatial resolution. Nondimensional color scale units (PV is normalized using the length scale 3×10^4 m, corresponding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

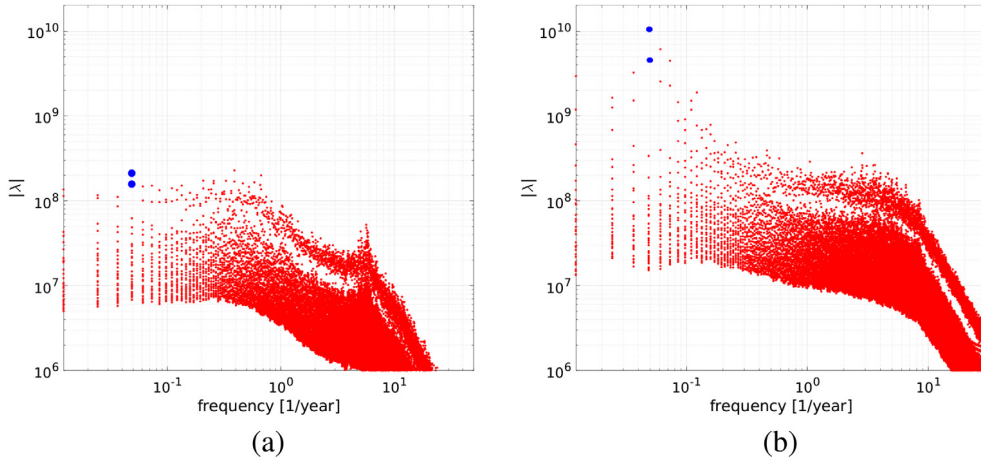


Fig. 2. Temporal spectral content of the reference solutions with: (a) 129^2 and (b) 513^2 grids. Shown are the 30 largest values of $|\lambda|$ per frequency, as given by the DAHD power spectrum of the upper-layer PV anomalies. The blue dots in panel (b) indicate maximum of the broadband spectral peak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

3. Scale decomposition of the high-resolution solution

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

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

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

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

where the overbar and prime indicate the large-scale and eddy components, respectively. With this in mind, we decomposed the high-resolution streamfunctions ψ_i by the moving-average square filter of size W ; and the corresponding PV anomalies are obtained by differentiation (akin Eq. (3)). We justify our choice of W by focusing on mesoscale eddies, which are scaled by the first baroclinic Rossby deformation radius, but we also admit that the problem contains many length scales and they vary geographically making the flow decomposition a difficult and open problem. The problem stems from the fact that for linear flows (when all the active scales are well separated in the Fourier spectra), the filter size should linearly depend on the ratio between the fine – and coarse – resolution grids. However, in our case, there is no separation between the active scales and the filter size is chosen based on the expected dynamical features we would like to filter out assuming the coarse-resolution model being unable to resolve them. In our case, these features are mesoscale eddies with length scales of order of the first baroclinic Rossby deformation radius ($\approx 10 - 100$ km).

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

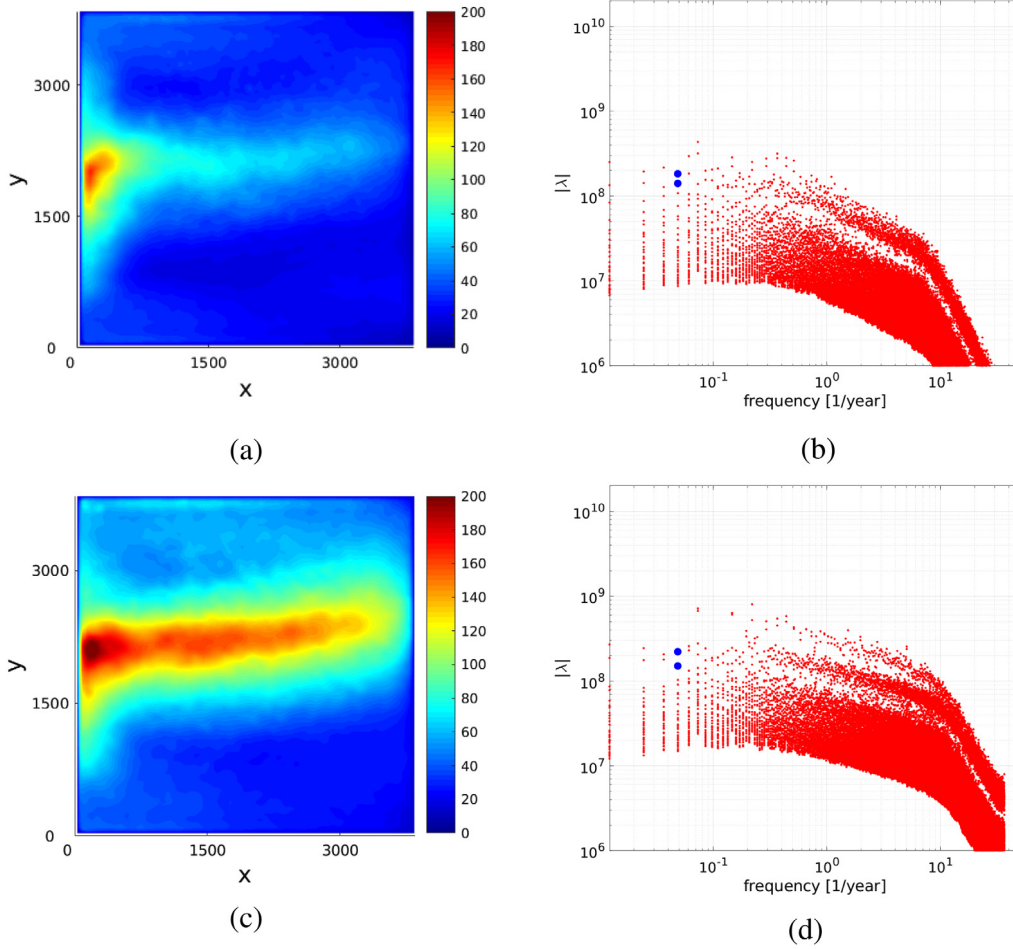


Fig. 3. Statistics of the upper-layer PV anomaly field for the low-resolution augmented solution (129^2 grid) obtained by feeding the true eddy field extracted with the $W = 21$ filter): (a) standard deviation showing partial reconstruction of the eastward jet extension; (b) temporal spectral content 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.

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

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

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

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

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

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

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

With only the eddies being fed to the augmented model, the external information is subtler, which makes it harder for the low-resolution

model to resolve desired dynamics resembling the fine-resolution reference solution such that the eastward extension of the jet is noticeably rectified and the low-frequency variability is present. At the same time, gauging the possibility of reducing the amount of data necessary for successful parameterization and errors introduced due to the incompleteness of the data is practically important.

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

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

where the small-scale (eddy) fields ψ'_i, q'_i are taken from the high-resolution data, and the prognostic low-resolution, large-scale variables ψ_i, q_i are continuously updated *online* during numerical integration of the model. We used all 90 years of the daily output to extract the eddy fields and then linearly interpolated them in time in-between the data records. An important issue of determining the minimal length of the eddy history for the quality augmentation of the low-resolution model is left outside the scope of the paper and will be addressed elsewhere.

We assessed the quality of the augmented low-resolution solution by looking into the simulated eastward jet region, focusing on its large-scale circulation patterns (evinced by the standard deviation in time) and LFV. The augmented-model eastward jet has improved but is still substantially different from the reference truth, as can be seen by comparing Figs. 3a and 1a. Similarly large discrepancies are seen in the augmented-model DAHD spectrum (Fig. 3b), which completely lacks the LFV. The interactive eddy forcing (Fig. 4a) can be significantly

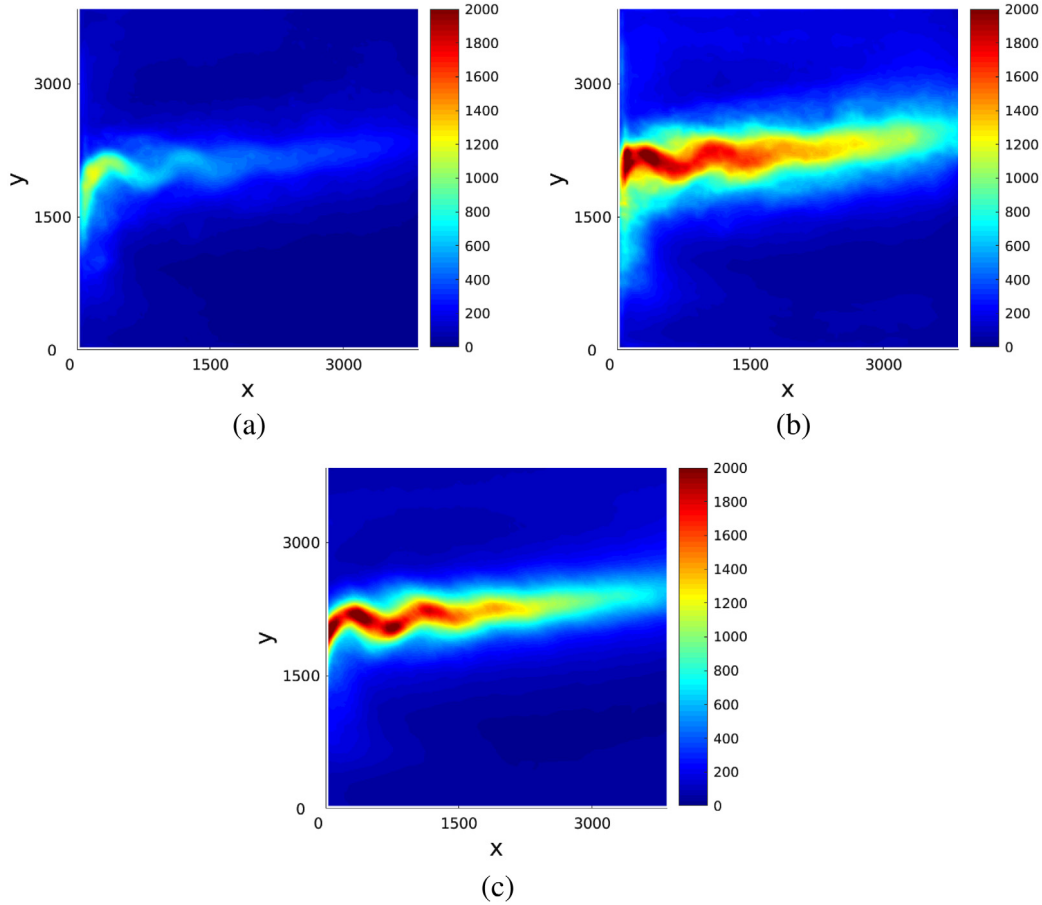


Fig. 4. Standard deviations of different eddy forcings: (a) on-line eddy forcing from the solution augmented with eddies extracted with filter size $W = 21$; (b) same as (a), but for $W = 41$; (c) true (offline) eddy forcing, as in Ryzhov et al. (2019)). The on-line eddy forcing in (a) is about 4 times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

less efficient because it is noticeably weaker than the true eddy forcing (Fig. 4c). We checked this by considering the more energetic eddy field extracted with the larger filter size $W = 41$ (Fig. 4b), but although the resulting eddy forcing is as intensive as the true one, the augmented model is still incapable of generating the LFV as implied by the DAHD spectrum (Fig. 3d). From this, we conclude that feeding even the most complete eddy fields into the model is still not sufficient for augmenting the solution. So, one has to use additional information from the high-resolution data to induce the LFV.

It has been already established (Ryzhov et al., 2019) that the true (off-line) eddy-forcing (Fig. 4b) generates the LFV in the augmented solution; therefore, we know that one way or another the model can be successfully augmented with the right amount of the extra information. One way to add this information is by interactively projecting the augmented solution onto the leading, true large-scale EOFs, and this can be viewed as a weak statistical constraint imposed by the filtering. The corresponding set of EOFs are obtained through the standard singular value decomposition, such that

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

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

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

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

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

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

The additionally filtered model solutions now exhibit the LFV as diagnosed by DAHD spectra shown in Fig. 5a for $W = 21$ and Fig. 5b for $W = 41$. It is worth noting that even in the solution augmented with weaker eddies ($W = 21$) the LFV is also reproduced, albeit it is not as energetic as with the stronger eddies ($W = 41$). The eastward jet extension is also reproduced similarly to the case without large-scale filtering (see Fig. 3).

In addition to the detailed DAHD spectral space–time diagnostic of PV anomaly field, it is also useful to consider the manifestation of LFV in the total potential energy, which is a global characteristic of the solution. Fig. 6 shows the Fourier spectral analysis of the potential energy time series by the standard Multitaper method (Percival and Walden, 1993), which reveals broadband LFV peaks at frequency

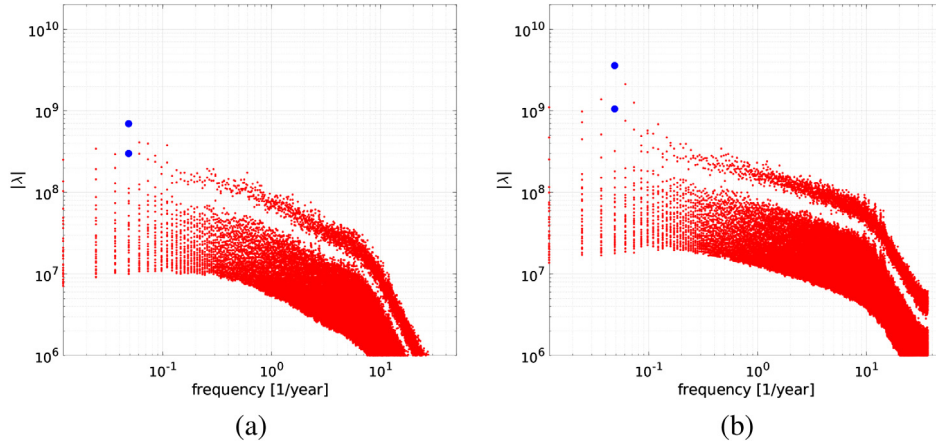


Fig. 5. The DAHD temporal spectra of the upper-layer PV anomaly field in augmented and additionally filtered model solutions: (a) $W = 21$ (weaker eddies) ; (b) $W = 41$ (stronger eddies). The LFV (see the peaks with the blue dots) is now present in both solutions, and it is more intensive with stronger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

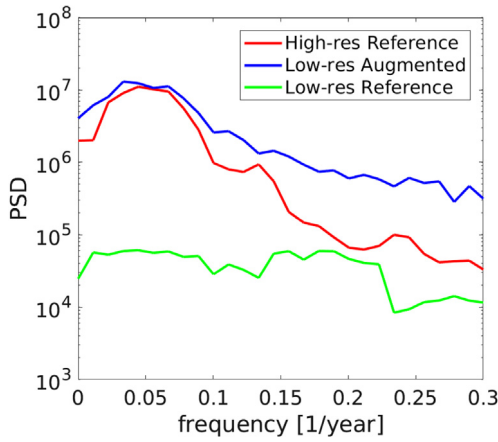


Fig. 6. Power spectrum density (PSD) of the potential energy by the Multitaper method, featuring the energetic and broadband LFV with the main period of ≈ 17 years, in both the reference high-resolution solution and augmented low-resolution solution (supplied by the eddy field obtained with filter $W = 41$ 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.

$\approx 0.06 \text{ year}^{-1}$ (about 17 years period), both for the reference high-resolution and augmented low-resolution solutions, whilst the reference low-resolution solution features no LFV with a mostly flat spectrum. Due to the projection, the augmented solution acquires oversaturated high frequencies near the LFV peak; this may be dealt with by carefully selecting the projection basis of the filtering procedure so to filter out spurious small-scale effects and is beyond the scope of the current study as we aimed at imbuing the coarse-resolution solution with the correct LFV.

Finally, we would like to emphasize that feeding the eddies to induce the augmenting eddy forcing in the low-resolution model (Eq. (9)) is absolutely necessary for generating the LFV, and we verified this by turning it off. If the filtering based on the EOF projection procedure is applied alone, it does not augment the solution thus confirming that the main component of the parameterization is the eddy forcing.

5. Statistical emulation of the eddy field

Here we developed data-driven statistical emulators of the true eddy field for feeding them into the low-resolution model instead of the original high-resolution eddy fields. The number of statistical emulation methods has recently surged, including stochastic approaches in climate science (Penland and Matrosova, 2001; Strounine et al.,

2010; Franzke et al., 2015; Kondrashov et al., 2015; Chen et al., 2016; Palmer, 2019; Seleznev et al., 2019; Foster et al., 2020), as well as other machine-learning (deep learning) methods developed for fluid dynamics applications (Brunton et al., 2020; Bolton and Zanna, 2019). The detailed analysis of emulated eddy fields is beyond the scope of this study, and in the context of assessing the skill of our emulators we focus solely on one of the central problems in climate ocean model simulations, namely, the correct rectification of the eddy field’s impact on the large-scale circulation. Thus we aimed for the solution of the low-resolution model, when augmented by an emulated eddy field, to be able to reproduce the long-term statistics of the high-resolution reference solution. We utilized the same skill measures as for the true eddy field explored in previous section. These are the geometrical shape of the large-scale circulation patterns, as well as the manifestation of the LFV.

We used a 30 000-day long high-resolution dataset of the eddy stream function $\tilde{\psi}$ for the three layers combined. The dataset is then coarse-grained onto the low spatial resolution (129×129), and further compressed by the PCA. We retained the leading 1000 PCs that account for $\approx 98\%$ of the variability.

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

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

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

The emulated PCs are obtained by initializing the model from the first data point of the training interval and by running it for 30 000 days. The eddy field is reconstructed in space from the emulated PCs by using the EOF basis, and then it is fed into the low resolution model in our augmentation procedure. While this basic emulator of the eddy field yields a fairly reasonable geometrical structure of the jet extension in the augmented solution (Fig. 7a), it does not induce the LFV as evident by the flat spectral density curve of the full potential energy (Fig. 7b), which is also similar to the non-augmented low-resolution solution.

A closer analysis shows that the lack of the LFV in the augmented solution is related to the spectral content of the emulated eddy field, in which energy at low frequencies is underestimated in comparison to the true eddy field. In turn, because the LFV in the true eddy

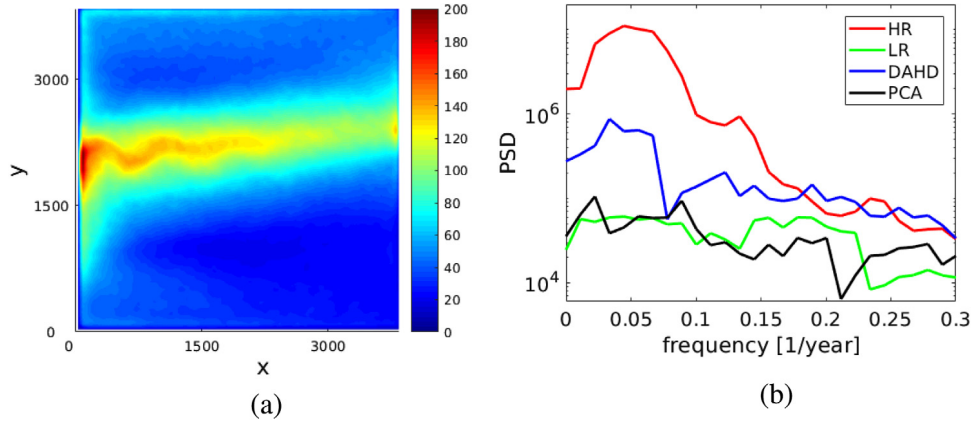


Fig. 7. (a) Standard deviation of the upper-layer PV anomalies in the augmented solution with an artificial eddy field emulated by a PCA-based linear model (Eq. (13)) (the periodical projection onto the 2000 EOFs of the large-scale “truth” basis is applied as well). The pattern of the standard deviation 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.

field is considerably weaker than in the true reference solution, it is challenging to capture it by an emulator based on PCA PCs, which typically mix different temporal scales.

The DAHD method (Section 2.2 and Appendix A) provides a novel emulation alternative, as it combines identification of frequency-ranked modes and their efficient modeling. It extracts pairs of data-adaptive harmonic modes (DAHMs) that form an orthonormal set of spatial patterns oscillating harmonically in time, and, thus, represent global monochromatic space–time filters. Projection of the dataset onto DAHMs yields pairs of narrowband time series of data-adaptive harmonic coefficients (DAHCs), which are modulated in amplitude, but do not mix temporal scales.

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

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

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 an additive noise. Furthermore, multiple SL-oscillators associated with the same non-zero frequency are linearly coupled and synchronized across frequencies by the pairwise-correlated white noise, while the model parameters are estimated by a regression with constraints (see Appendix B for numerical details). The original dataset with its multiple time scales can be modeled in a computationally efficient manner since the contribution of each temporal frequency is simulated in parallel.

Kondrashov et al. (2018a) developed a stochastic DAHD emulator for the LFV in the model considered, and here we extended these results to the eddies. We used the leading $d = 100$ PCs of the eddy streamfunction capturing $\approx 70\%$ of the variance and applied the DAHD with the embedding window of $M = 100$ days. Then, we fit the model of coupled $d = 100$ stochastic oscillators for the DAHCs and obtained their emulations for the $M = 100$ frequencies. After emulated DAHCs were back-transformed into the space–time eddy field by using DAHMs and EOFs, and combined across all the emulated frequencies, we fed the outcome into the augmented model. The geometrical shape of the augmented solution is again reproduced fairly well, and it is very similar to Fig. 7a (not shown for brevity). Furthermore, since the LFV is now better captured in the emulated eddy field (compare to the high-resolution “truth”), it is also induced in the augmented solution (Fig. 7b), albeit it is less energetic than when the true eddy field is used (see Fig. 6).

6. Conclusions

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

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

In the high-resolution reference solution both of the key features are well represented, whereas the low-resolution reference solution lacks any of them. Motivation for including (1) is straightforward, because any eddy parameterization is, first of all, tested for its ability to simulate the large-scale climatological fields. Motivation for including (2) is to test the ability of the parameterization to simulate intrinsic climate variabilities similar to the relatively well understood interdecadal variability featured in our model. Our hope is that testing mesoscale eddy parameterization skills will eventually include climate variability signals as the standard test beds.

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

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

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

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

This approach worked well, and we experimentally found the optimal number of the EOFs and the optimal frequency of the applied filtering procedure, so that the LFV was almost fully recovered. Since the filtering can be applied infrequently (about every 100 days in our case) rather than continuously, which is also possible, its computational cost is nearly negligible. However, the exact amount of information needed from the high-resolution “truth” for a correct rectification of the LFV remains unknown and its assessment should be addressed elsewhere. We hypothesized that this information should contain correct correlations between the eddy and large-scale fields. We also demonstrated that the filtering was of secondary importance relative to the supplied eddy forcing, because when the latter was switched off, the filtering alone was not capable of augmenting the solution to any acceptable level.

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

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

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

We thank Dr J. Maddison and one anonymous reviewer for constructive comments that helped improve this manuscript. This research was supported by the National Science Foundation (NSF), USA grants OCE – 1658357 and the Natural Environment Research Council (NERC) grant NE/R011567/1. Pavel Berloff also gratefully acknowledges funding by NERC, UK Grant No. NE/T002220/1 and Leverhulme Trust, UK Grant No. RPG – 2019 – 024 and the Moscow Center for Fundamental and Applied Mathematics (supported by the Agreement 075 – 15 – 2019 – 1624 with the Ministry of Education and Science of the Russian Federation). DAHD analysis was supported by the Russian Science Foundation (Grant No. 18 – 12 – 00231). We would like to acknowledge the high-performance computing support from Cheyenne (doi: 10.5065/D6RX99HX) provided by NCAR’s Computational and Information Systems Laboratory, sponsored by the NSF, USA. The DAHD Toolbox is available at: <http://research.atmos.ucla.edu/tcd/dkondras/Software.html>

Appendix A. Data-adaptive harmonic decomposition (DAHD)

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

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

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

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

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

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

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

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

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

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

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

Appendix B. Frequency-ranked stochastic emulators

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

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

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

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

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

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

References

Agarwal, N., Ryzhov, E., Kondrashov, D., Berloff, P., 2020. Scale-aware flow decomposition and statistical analysis of the eddy forcing. submitted for publication.

Arnold, H.M., Moroz, I.M., Palmer, T.N., 2013. Stochastic parametrizations and model uncertainty in the Lorenz 96 system. *Phil. Trans. R. Soc. A* 371 (1991), 20110479.

Bachman, S.D., Fox-Kemper, B., Pearson, B., 2017. A scale-aware subgrid model for quasi-geostrophic turbulence. *J. Geophys. Res.: Oceans* 122 (2), 1529–1554.

Berloff, P., 2005. On dynamically consistent eddy fluxes. *Dyn. Atmos. Oceans* 38, 123–146.

Berloff, P., 2015. Dynamically consistent parameterization of mesoscale eddies. Part I: simple model. *Ocean Model.* 87, 1–19.

Berloff, P., 2016. Dynamically consistent parameterization of mesoscale eddies. Part II: eddy fluxes and diffusivity from transient impulses. *Fluids* 1 (3A).

Berloff, P., 2018. Dynamically consistent parameterization of mesoscale eddies. Part III: Deterministic approach. *Ocean Model.* 127, 1–15.

Berloff, P., Hogg, A., Dewar, W., 2007. The turbulent oscillator: A mechanism of low-frequency variability of the wind-driven ocean gyres. *J. Phys. Oceanogr.* 37, 2363–2386.

Berloff, P.S., McWilliams, J., 1999. Large-scale, low-frequency variability in wind-driven ocean gyres. *J. Phys. Oceanogr.* 29, 1925–1949.

Bolton, T., Zanna, L., 2019. Applications of deep learning to ocean data inference and subgrid parameterization. *J. Adv. Modelling Earth Syst.* 11 (1), 376–399.

Brunton, S.L., Noack, B.R., Koumoutsakos, P., 2020. Machine learning for fluid mechanics. *Annu. Rev. Fluid Mech.* 52 (1).

Chekroun, M.D., Kondrashov, D., 2017. Data-adaptive harmonic spectra and multilayer Stuart–Landau models. *Chaos* 27, 093110.

Chen, C., Cane, M.A., Henderson, N., Lee, D.E., Chapman, D., Kondrashov, D., Chekroun, M.D., 2016. Diversity, nonlinearity, seasonality, and memory effect in ENSO simulation and prediction using empirical model reduction. *J. Clim.* 29 (5), 1809–1830.

Chorin, A.J., Lu, F., 2015. Discrete approach to stochastic parametrization and dimension reduction in nonlinear dynamics. *PNAS* 112 (32), 9804–9809.

Crommelin, D., Vanden-Eijnden, E., 2008. Subgrid-scale parameterization with conditional Markov chains. *J. Atmos. Sci.* 65 (8), 2661–2675.

Dijkstra, H.A., 2013. *Nonlinear Climate Dynamics*. Cambridge University Press, Cambridge, UK.

Dijkstra, H.A., 2018. A normal mode perspective of intrinsic ocean-climate variability. *Annu. Rev. Fluid Mech.* 48, 341–363.

Fatkullin, I., Vanden-Eijnden, E., 2004. A computational strategy for multiscale systems with applications to Lorenz 96 model. *J. Comput. Phys.* 200 (2), 605–638.

Foster, D., Comeau, D., Urban, N.M., 2020. A Bayesian approach to regional decadal predictability: Sparse parameter estimation in high-dimensional linear inverse models of high-latitude sea surface temperature variability. *J. Clim.* 33 (14), 6065–6081.

Franzke, C.L.E., O’Kane, T.J., Berner, J., Williams, P.D., Lucarini, V., 2015. Stochastic climate theory and modeling. *Wiley Interdiscip. Rev. Clim. Change* 6 (1), 63–78.

Frederiksen, J.S., 1999. Subgrid-scale parameterizations of eddy-topographic force, eddy viscosity, and stochastic backscatter for flow over topography. *J. Atmos. Sci.* 56 (11), 1481–1494.

Frederiksen, J.S., O’Kane, T.J., Zidikheri, M.J., 2012. Stochastic subgrid parameterizations for atmospheric and oceanic flows. *Phys. Scr.* 85 (6), 068202.

Gent, P.R., McWilliams, J.C., 1990. Isopycnal mixing in ocean circulation models. *J. Phys. Oceanogr.* 20 (1), 150–155.

Hasselmann, K., 1988. PIPs and POPs: The reduction of complex dynamical systems using principal interaction and oscillation patterns. *J. Geophys. Res.: Atmos.* 93 (D9), 11015–11021.

Jansen, M.F., Held, I.M., 2014. Parameterizing subgrid-scale eddy effects using energetically consistent backscatter. *Ocean Model.* 80, 36–48.

Jansen, M.F., Held, I.M., Adcroft, A., Hallberg, R., 2015. Energy budget-based backscatter in an eddy permitting primitive equation model. *Ocean Model.* 94, 15–26.

Karabasov, S., Berloff, P., Goloviznin, V., 2009. CABARET in the ocean gyres. *Ocean Model.* 30 (2), 155–168.

Kondrashov, D., Berloff, P., 2015. Stochastic modeling of decadal variability in ocean gyres. *Geophys. Res. Lett.* 42, 1543–1553.

Kondrashov, D., Chekroun, M.D., 2018. Data-adaptive harmonic analysis and modeling of solar wind-magnetosphere coupling. *J. Atmos. Sol.-Terr. Phys.*

Kondrashov, D., Chekroun, M.D., Berloff, P., 2018a. Multiscale Stuart–Landau emulators: Application to wind-driven ocean gyres. *Fluids* 3, 21.

Kondrashov, D., Chekroun, M.D., Ghil, M., 2015. Data-driven non-Markovian closure models. *Physica D* 297, 33–55.

Kondrashov, D., Chekroun, M.D., Ghil, M., 2018b. Data-adaptive harmonic decomposition and prediction of Arctic sea ice extent. *Dyn. Stat. Clim. Syst.* 3 (1).

Kondrashov, D., Chekroun, M.D., Yuan, X., Ghil, M., 2018c. Data-adaptive Harmonic Decomposition and Stochastic Modeling of arctic Sea Ice. In: Tsonis, A.A. (Ed.), *Advances in Nonlinear Geosciences*. Springer International Publishing, Cham, pp. 179–205.

Kondrashov, D., Kravtsov, S., Robertson, A.W., Ghil, M., 2005. A hierarchy of data-based ENSO models. *J. Clim.* 18 (21), 4425–4444.

Kondrashov, D., Ryzhov, E.A., Berloff, P., 2020. Data-adaptive harmonic analysis of oceanic waves and turbulent flows. *Chaos* 30 (6), 061105.

Kravtsov, S., Berloff, P., Dewar, W., Ghil, M., McWilliams, J., 2006. Dynamical origin of low-frequency variability in a highly nonlinear midlatitude coupled model. *J. Clim.* 19, 6391–6408.

Kravtsov, S., Kondrashov, D., Ghil, M., 2005. Multi-level regression modeling of nonlinear processes: Derivation and applications to climatic variability. *J. Clim.* 18, 4404–4424.

Li, H., von Storch, J.-S., 2013. On the fluctuating buoyancy fluxes simulated in a OGCM. *J. Phys. Oceanogr.* 43 (7), 1270–1287.

Maddison, J.R., Marshall, D.P., Shipton, J., 2015. On the dynamical influence of ocean eddy potential vorticity fluxes. *Ocean Model.* 92, 169–182.

Majda, A.J., Timofeyev, I., Vanden Eijnden, E., 1999. Models for stochastic climate prediction. *Proc. Natl. Acad. Sci. USA* 96 (26), 14687–14691.

- Mak, J., Maddison, J.R., Marshall, D.P., Munday, D.R., 2018. Implementation of a geometrically informed and energetically constrained mesoscale eddy parameterization in an ocean circulation model. *J. Phys. Oceanogr.* 48 (10), 2363–2382.
- Marshall, D.P., Maddison, J.R., Berloff, P.S., 2012. A framework for parameterizing eddy potential vorticity fluxes. *J. Phys. Oceanogr.* 42 (4), 539–557.
- Palmer, T.N., 2019. Stochastic weather and climate models. *Nat. Rev. Phys.* 1 (7), 463–471.
- Penland, C., Matrosova, L., 2001. Expected and actual errors of linear inverse model forecasts. *Mon. Weather Rev.* 129 (7), 1740–1745.
- Percival, D.B., Walden, A.T., 1993. *Spectral Analysis for Physical Applications*. Cambridge University Press.
- Porta Mana, P., Zanna, L., 2014. Toward a stochastic parametrization of ocean mesoscale eddies. *Ocean Model.* 79, 1–20.
- Preisendorfer, R.W., 1988. *Principal Component Analysis in Meteorology and Oceanography*. Elsevier, New York, p. 425.
- Ryzhov, E., Kondrashov, D., Agarwal, N., Berloff, P., 2019. On data-driven augmentation of low-resolution ocean model dynamics. *Ocean Model.* 142, 101464.
- Schmid, P.J., 2010. Dynamic mode decomposition of numerical and experimental data. *J. Fluid Mech.* 656, 5–28.
- Seleznev, A., Mukhin, D., Gavrilov, A., Loskutov, E., Feigin, A., 2019. Bayesian framework for simulation of dynamical systems from multidimensional data using recurrent neural network. *Chaos* 29 (12), 123115.
- Shevchenko, I.V., Berloff, P.S., 2015. Multi-layer quasi-geostrophic ocean dynamics in eddy-resolving regimes. *Ocean Model.* 94, 1–14.
- Shevchenko, I., Berloff, P., 2016. Eddy backscatter and counter-rotating gyre anomalies of midlatitude ocean dynamics. *Fluids* 1 (3).
- Shevchenko, I., Berloff, P., Guerrero-Lopez, D., Roman, J., 2016. On low-frequency variability of the midlatitude ocean gyres. *J. Fluid Mech.* 795, 423–442.
- Strounine, K., Kravtsov, S., Kondrashov, D., Ghil, M., 2010. Reduced models of atmospheric low-frequency variability: Parameter estimation and comparative performance. *Physica D* 239 (3), 145–166.
- Viebahn, J., Crommelin, D., Dijkstra, H., 2019. Toward a turbulence closure based on energy modes. *J. Phys. Oceanogr.* 49 (4), 1075–1097.
- von Storch, H., Bürger, G., Schnur, R., von Storch, J.-S., 1995. Principal oscillation patterns: A review. *J. Clim.* 8 (3), 377–400.
- Ying, Y.K., Maddison, J.R., Vanneste, J., 2019. Bayesian inference of ocean diffusivity from Lagrangian trajectory data. *Ocean Model.* 140, 101401.
- Zanna, L., Porta Mana, P., Anstey, J., David, T., Bolton, T., 2017. Scale-aware deterministic and stochastic parametrizations of eddy-mean flow interaction. *Ocean Model.* 111, 66–80.