

Earth and Space Science



RESEARCH ARTICLE

10.1029/2018EA000519

Key Points:

- Machine learning is applied to create a geography of global dynamical regions, emergent dynamics are discussed
- A barotropic vorticity framework is sufficient for 93% of the ocean, with nonlinear terms have a small extent but could impact circulation
- Robust and novel application of machine learning techniques are demonstrated to offer insight into large-scale ocean physical regimes

Correspondence to:

M. Sonnewald,
maike_s@mit.edu

Citation:

Sonnewald, M., Wunsch, C., & Heimbach, P. (2019). Unsupervised learning reveals geography of global ocean dynamical regions. *Earth and Space Science*, 6, 784–794. <https://doi.org/10.1029/2018EA000519>

Received 14 NOV 2018

Accepted 24 FEB 2019

Accepted article online 13 MAR 2019

Published online 21 MAY 2019

Author Contributions

Conceptualization: Maïke Sonnewald, Carl Wunsch

Funding Acquisition: Patrick Heimbach

Methodology: Maïke Sonnewald

Writing - Original Draft: Maïke Sonnewald

Formal Analysis: Maïke Sonnewald

Investigation: Maïke Sonnewald

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Writing - review & editing: Maïke Sonnewald, Carl Wunsch, Patrick Heimbach

Resources: Carl Wunsch, Patrick Heimbach

Unsupervised Learning Reveals Geography of Global Ocean Dynamical Regions

Maïke Sonnewald^{1,2} , Carl Wunsch^{1,2}, and Patrick Heimbach³

¹Department of Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA, ²Department of Earth and Planetary Sciences, Harvard University, Cambridge, MA, USA, ³Oden Institute for Computational Engineering and Sciences, University of Texas at Austin, Austin, TX, USA

Abstract Dynamically similar regions of the global ocean are identified using a barotropic vorticity (BV) framework from a 20-year mean of the Estimating the Circulation and Climate of the Ocean state estimate at 1° resolution. An unsupervised machine learning algorithm, *K*-means, objectively clusters the standardized BV equation, identifying five unambiguous regimes. Cluster 1 covers $43 \pm 3.3\%$ of the ocean area. Surface and bottom stress torque are balanced by the bottom pressure torque and the nonlinear torque. Cluster 2 covers $24.8 \pm 1.2\%$, where the beta effect balances the bottom pressure torque. Cluster 3 covers $14.6 \pm 1.0\%$, characterized by a “Quasi-Sverdrupian” regime where the beta effect is balanced by the wind and bottom stress term. The small region of Cluster 4 has baroclinic dynamics covering $6.9 \pm 2.9\%$ of the ocean. Cluster 5 occurs primarily in the Southern Ocean. Residual “dominantly nonlinear” regions highlight where the BV approach is inadequate, found in areas of rough topography in the Southern Ocean and along western boundaries.

Plain Language Summary A geography of the global ocean is presented of dynamical regions found using unsupervised machine learning techniques. A bottom-up approach is used to identify emergent patterns in the modern global ocean. Their existence demonstrates commonalities that lead to methods for understanding the global circulation. Five areas vary from a depth coherent flow, quasi-Sverdrupian, interior, interior with a strong vertical component, and an interior flow specific to the Southern Ocean. In some regions nonlinear terms are important, which will be the subject of future work at higher resolution.

1. Motivation

Before the advent of modern observational and modeling techniques, understanding of the physical/dynamical state of the ocean focused on large-scale quasi-laminar descriptions such as Sverdrup balance, abyssal recipes, or Stommel-Arons flows (Munk, 1950; Munk & Palmén, 1940; Stommel, 1948). Recent advances in instrumentation and modeling capability have revealed that ocean physics is characterized by complex spatial and temporal variability. It is possible that every location in the ocean has a unique physical state depending upon many factors, including local topography, meteorology, proximity to eastern and western boundaries, or latitude, rendering any global scale interpretation lacking in general applicability.

The ocean is spatially and temporally diverse, but delineating spatial and temporal commonalities and continuities are central to understanding emergent patterns that lead to a global geography of dynamical regimes. Dominant physics underlying the emergent patterns become evident when common features are identified. This note's purpose is to explore unsupervised machine learning as a method for depicting and understanding the gross features of global oceanic physics. The study is restricted to the barotropic vorticity (BV) balance of a time mean global circulation as calculated from a noneddy resolving state estimate (Fukumori et al., 2018). Our approach appears to be both interesting and useful and is readily generalized to far more complex oceanic states. In unsupervised machine learning approaches, data are not pre-labeled or pre-grouped (Bishop, 2006). The noneddy resolving case is explored in this initial work, analogous to Coupled Model Intercomparison Project Phases 5 and 6 (CMIP5 and CMIP6) efforts (Church & Miles, 2013; Eyring et al., 2016; Stouffer et al., 2017). The presented analysis is intended to provide a description relevant to the CMIP climate models.

©2019. The Authors.

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

The presence of differing dynamical regimes is already suggested by the known structures of the wind-stress forcing and the geometry of the ocean basins, including the underlying topography. Classifying and identifying regions in the world ocean is done here using the variations in the dominant terms of the BV budget, following the demonstration that the global budget can be closed. Schoonover et al. (2016) and Yeager (2015) have assessed the dynamics of the BV budget focusing on the North Atlantic Ocean variability and sensitivity to resolution, respectively. Here the procedures are global.

Distinct geographical physics were demonstrated by, for example, Xu and Fu (2012), who used altimeter data to show differing spatial regimes of geostrophic turbulence. Their global patterns are presumably connected to circulation structures, planetary waves, and topography. Similarly, Hughes and Williams (2010) and Sonnewald et al. (2018) used linear statistical models to infer patterns suggestive of global regimes. The present work extends the pattern determination methodology, as working only with data from the surface limits the application to a comprehensive assessment of global dynamical regimes. Use of the Estimating the Circulation and Climate of the Ocean (ECCO) state estimate extends the previous surface data analyses to the full three dimensions of the ocean.

Objective classification via K -means clustering allows unbiased identification of patterns in data. This form of unsupervised machine learning is common in many fields, ranging from pharmaceuticals to engineering (Breuhl et al., 1999; Hauser & Rybakowski, 1997; Kulis & Jordan, 2012). Ardyna et al. (2017) applied a similar method to identify regions with distinct biological activity, and K -means have been used to identify key regions for data collection to build maps of nitrate in the Southern Ocean (Y.-C. Liang, personal communication, February 15, 2018). Applications in the earth sciences have been explored both in the prognostic and diagnostic sense (Krasnopolsky et al., 2013; Schneider et al., 2017).

2. Methods

2.1. The ECCO Version 4 State Estimate

The BV equation is applied to version 4, release 2, of the ECCO (ECCOv4) global state estimate, described by Forget et al. (2015) and Wunsch and Heimbach (2013), see also ECCO Consortium (2017a, 2017b). The state estimate has a nominally 1° resolution, available at ecco.jpl.nasa.gov. A least squares with Lagrange multipliers (4DVAR) approach is used to obtain the state estimate. The result is a *free-running* version of the MIT General Circulation Model (MITgcm; Adcroft et al., 2004), with adjusted initial and boundary conditions and internal model parameters. The ECCO state satisfies basic conservation laws for enthalpy, salt, volume, and momentum while remaining largely within error estimates of a diverse set of global data (Stammer et al., 2016; Wunsch & Heimbach, 2007, 2013). Regions without data are brought into consistency in a dynamically consistent way using the dynamics, still relying on parameterizations but avoiding the use of untested statistical hypotheses or infilling e.g., Reynolds et al. (2013).

2.2. Barotropic Vorticity

The momentum and continuity equations of an ocean in a thin shell on a rotating sphere are

$$\partial_t \mathbf{u} + f \mathbf{k} \times \mathbf{u} = -\frac{1}{\rho_0} \nabla_h p + \frac{1}{\rho_0} \partial_z \tau + \mathbf{a} + \mathbf{b}, \quad \partial_z p = -g\rho, \quad (1)$$

$$\nabla_h \cdot \mathbf{u} + \partial_z w = 0 \quad (2)$$

Pressure, gravity, density, and vertical shear stress are denoted p , g , ρ , and τ , respectively, with ρ_0 the reference density; the three-dimensional velocity field $\mathbf{v} = (\mathbf{u}, \mathbf{v}, \mathbf{w}) = (\mathbf{u}, \mathbf{w})$; the gradient $\nabla = (\nabla_h, \partial_z)$; the unit vector is denoted \mathbf{k} ; planetary vorticity is a function of latitude ϕ in $f \mathbf{k} = (0, 0, 2\Omega \sin \phi)$; the viscous forcing by vertical shear is denoted $\partial_z \tau$; the nonlinear torque is \mathbf{a} , and the horizontal viscous forcing \mathbf{b} includes subgrid-scale parameterizations. Assuming a steady state, the vertical integral from the surface $z = \eta(x, y, t)$ to the water depth below the surface $z = H(x, y)$ is

$$\beta V = \frac{1}{\rho_0} \nabla p_b \times \nabla H + \frac{1}{\rho_0} \nabla \times \tau + \nabla \times \mathbf{A} + \nabla \times \mathbf{B}, \quad (3)$$

where $\nabla \cdot \mathbf{U} = 0$, $\mathbf{U} \cdot \nabla f = \beta V$, the bottom pressure is denoted p_b , $\mathbf{A} = \int_H^\eta \mathbf{a} dz$, and $\mathbf{B} = \int_H^\eta \mathbf{b} dz$. The curl operator $\nabla \times$ yields a scalar, representing the vertical component of the operator. The left-hand side of equation (3) is the planetary vorticity advection term, while the right-hand side of equation (3) is the bottom pressure torque (BPT), the wind and bottom stress curl, the nonlinear torque, and the viscous torque,

Table 1*Percentage of Area Covered by the Area-Specific Balance of the BV Equation (3) and the Corresponding Map Figure*

Cluster	Area	Leading terms
1	43 ± 3.3%, Depth coherent (Figure 3a)	$\nabla \times \tau_{sb} + \nabla \times \mathbf{A} \approx \nabla p_b \times \nabla H$ (Figure 3b)
2	24.8 ± 1.2%, Interior flow (Figure 3c)	$\nabla \times \tau_{sb} \approx \nabla p_b \times \nabla H + \nabla \cdot (f\mathbf{U})$ (Figure 3d)
3	14.6 ± 1%, Quasi-Sverdrupian (Figure 3e)	$\nabla \times \tau_{sb} \approx \nabla \cdot (f\mathbf{U})$ (Figure 3f)
4	6.9 ± 2.9%, Interior flow, vertical (Figure 4a)	$\nabla \times \tau_{sb} \approx \nabla \cdot (f\mathbf{U}) + \nabla p_b \times \nabla H$ (Figure 4b)
5	1.9 ± 1%, Interior flow, Southern Ocean (Figure 4c)	$\nabla \times \tau_{sb} \approx \nabla \cdot (f\mathbf{U}) + \nabla p_b \times \nabla H$ (Figure 4d)
6–50	8.9 ± 0.3%, Dominantly nonlinear (Figure 4e)	$\nabla \cdot (f\mathbf{U}) \approx \nabla \times \mathbf{A} + \nabla \times \tau_{sb}$ (Figure 4f)

Note. Leading order terms are sorted by magnitude, colors indicating if barotropic vorticity is added (red in font) or removed (blue in font) by the leading order term, the corresponding bar chart figure shows the full breakdown. The quoted percentage coverage and StD is the mean of 100 runs of the algorithm.

respectively. The subgrid-scale parameterization introduces a torque, which is included in the viscous torque term. Nonlinear torque is composed of three terms:

$$\nabla \times \mathbf{A} = \nabla \times \left[\int_{-H}^{\eta} \nabla \cdot (\mathbf{u}\mathbf{u}) dz \right] + [w\zeta]_{z=H}^{z=\eta} + [\nabla w \times \mathbf{u}]_{z=H}^{z=\eta}, \quad (4)$$

where $\mathbf{u}\mathbf{u}$ is a second-order tensor. The right-hand side of equation (4) represents the curl of the vertically integrated momentum flux divergence, the nonlinear contribution to vortex tube stretching, and the conversion of vertical shear to barotropic vorticity. Horizontal viscous forcing includes that induced by subgrid-scale parameterizations. Twenty-year averaged fields of the BV equation are used after a Laplacian smoother is applied, with an effective averaging range of three grid cells.

2.3. Unsupervised Learning: K-Means Clustering

Our goal is to determine the spatial patterns that correspond to various balances between the dominant terms of the BV equation. Various combinations of terms dominate the BV equation in different regions, and spatial patterns emerge, which we seek to determine. Clustering identifies groups in data based on how the data are distributed in parameter space, the dimensions of which are defined by a set of chosen variables (or “features” as they are also called in the machine learning literature such as ; Kubat, 2015). In this application, the terms of the BV equation are used to define the dimensions of the parameter space in which the clusters are identified. If groups of dominant terms are present that differ significantly, a robust separation into distinct “clusters” is feasible. For each cluster, the area-weighted histogram is calculated of the values of each term in the BV equation. Next, the clusters are sorted by the amount of geographical area that they cover using a sorting algorithm, associating the name of the cluster as listed in Table 1. This focuses on those clusters that featured clearly different balances between terms in the BV equation.

The terms in the BV equation were normalized so that each term individually has a zero mean and unit variance globally. This scaling of the variables ensures that the patterns in the variance of the variables form gridpoint to gridpoint are what is highlighted, rather than the relative magnitudes (Not applying normalization and feature scaling produced a poor result where clusters are harder to robustly identify.). Using the normalized and scaled fields, the clusters are more easily identified. The *K*-means algorithm (MacQueen, 1967) involves an iterative minimization of the sum of squares of the Euclidean distance partitioning of the hyperspace given by the terms in the BV equation:

$$J = \sum_{j=1}^K \sum_{i=1}^n \|\mathbf{x}_i^j - \mathbf{c}_j\|^2 \quad (5)$$

the number of *K* clusters is a free parameter that is chosen a priori; the cluster centers have random initial values scattered throughout the parameter space. The parameter \mathbf{x} is a vector field that is defined at every grid cell on a discretized sphere, with each element \mathbf{x}_i representing a five-dimensional vector on the model's horizontal grid, such that index *i* uniquely identifies a grid point on the sphere, with (lon,lat) = (ϕ_i, θ_i). The components (features) of each vector \mathbf{x}_i correspond to the five terms in the BV budget. Each cluster $j = 1, \dots, K$ is represented by the five-dimensional characterizing vector \mathbf{c}_j , and the *K*-means classification attributes each vector \mathbf{x}_i to a unique cluster \mathbf{c}_j , thus $\mathbf{x}_i = \mathbf{x}_i^j$. The distance between a data point \mathbf{x}_i^j and \mathbf{c}_j is given by $\|\mathbf{x}_i^j - \mathbf{c}_j\|^2$. Each data point is associated with the closest *K*-cluster, the position of \mathbf{c}_j is recalculated,

and the association reassessed until the solution converges. Assumptions regarding the covariance of the data are discussed in the appendix.

The solution is sensitive to the initialization and choice of K , and the algorithm partitions the parameter subspace using linear hyperplanes. This linearity constraint means that higher numbers of K can both assist in partitioning the subspace more appropriately, and isolate noise. The appendix demonstrates the small sensitivity of the result to the algorithm's initial random seed, and the impact of varying K . An optimal value of K is determined as $K > 35$ using the Akaike and Bayesian Information Criteria (AIC and BIC; Akaike, 1973). Information criteria provide a measure of the quality of a statistical model, rewarding increased likelihood across a data set and penalizing overfitting. AIC and BIC indicate robust regimes as they both asymptote in the bottom left panel of Figure 2, suggesting that no information is gained by further increasing K . $K = 50$ is used for the remaining analysis, where five clusters are individually analyzed as they are considered to represent somewhat classical dynamical regimes, and the remaining 45 clusters are taken together as a single “nonlinear regime.”

3. Results

The closure in ECCOv4 for the 20-year average of the BV terms in equation (3) is illustrated in Figure 1. Individual terms are of order $\pm 10^{-9} \text{ m s}^{-1}$ and the residual has magnitude of less than $\pm 10^{-12} \text{ m/s}$. For 36% of the ocean the residual is $\ll 1\%$ (Figure 1). This small residual permits going forward with confidence. Some numerical issues do exist on the continental shelf and in shallow water generally, but these regions only amount to 3% of the area of the global ocean and will be ignored.

Figure 1b illustrates where the beta term is important from equation (3). This term is balanced by the BPT term shown in Figure 1c and the wind and bottom stress BV terms shown in Figure 1d. The remainder is largely found in the nonlinear BV contributions seen in Figure 1e, with the lateral viscous dissipation largely being an order of magnitude smaller, apart from localized regions in the Southern Ocean. Wind and bottom stress BV terms in Figure 1d are largely zonally symmetric, with large patterns of negative BV in the Southern Ocean, and large gyre patterns visible in the Pacific and Atlantic basins. BPT in Figure 1c is associated with interactions with steep bathymetry. For example, in the Southern Ocean a large positive patch leads toward the Antarctic-Pacific ridge, with a negative patch beyond. This structure is consistent with vortex stretching as circulation crosses the ridge. Along Western Boundaries, BPT is positive to the west and negative just adjacent to the east, consistent with studies such as Myers et al. (1996). The BV of the nonlinear torque shown in Figure 1e is concentrated along the western edge of basins where WBCs are found, but it is less spatially coherent than the BPT term. Large activity stands out in the Southern Ocean region, particularly in the Atlantic sector. Lateral viscous dissipation is small.

Picking out globally coherent dynamical regimes, the K -means algorithm results are presented in Figure 2a, where the numbering on the color bar is arbitrary. The structure is mainly found in five named regimes, summarized in Table 1, which are briefly examined. Each is numbered, and a partially descriptive label is attached. Figure 2c shows the area and 2σ uncertainty covered by the named clusters.

Figures 3 and 4 isolate the geographical area (left column) for each cluster determined as distinct by the K -means algorithm. From this geographical area, the area-weighted histogram is computed (right column) across all terms of the BV equation. These histograms relate the clusters to different dynamical regimes, illustrated in Figure 1. Area averaging was done for comparison.

The largest cluster, accounting for 43% of the global ocean (Cluster 1 in Figure 3a), is determined by a balance between wind stress curl and bottom torque. This depth-coherent “negative wind curl/bottom torque” region is found primarily in zonal streaks in the tropics, and in a thin ribbon in the Southern Ocean mainly in the Pacific sector. In the Northern Hemisphere, Cluster 1 areas surround the subtropical and subpolar gyres. Large areas of the Arctic Seas are also in this Cluster. Figure 3b demonstrates that the balance of terms is dominantly between the input of negative vorticity by the wind-stress curl largely balanced by the positive input by the bottom interaction terms.

The next largest dynamical region covers 25% of the ocean area (Cluster 2, Figure 3c: Interior flow “positive wind curl/beta and bottom torque”), where the wind stress curl inputs positive vorticity, nearly balanced by

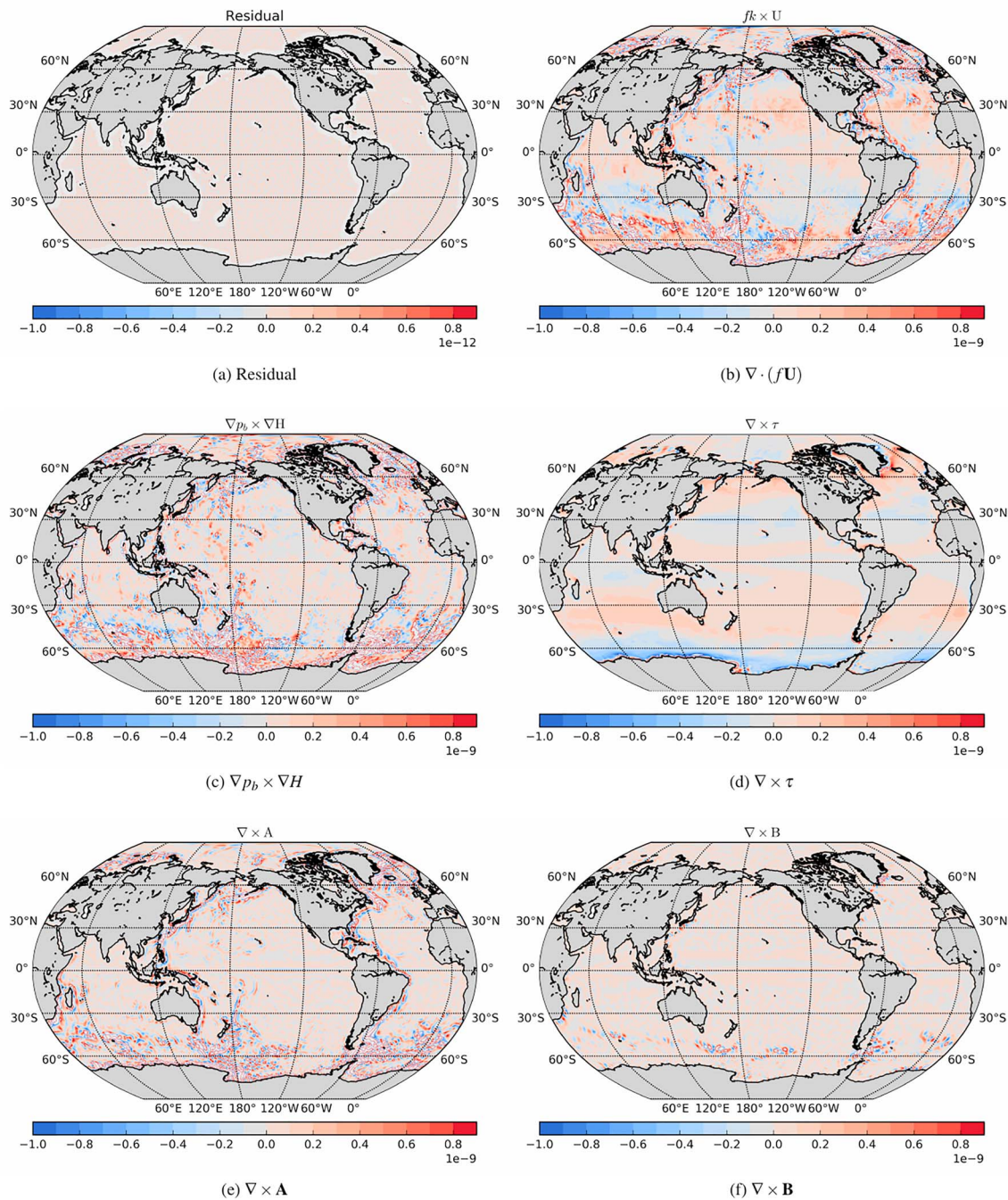


Figure 1. The breakdown of the barotropic vorticity budget (m/s) over 1992–2013 in the ECCOV4 State Estimate. From equation (3), Figure 1a is the planetary vorticity advection term, Figure 1b is the bottom pressure torque (BPT), Figure 1c is the wind and bottom stress curl, Figure 1d is the nonlinear torque, and the viscous torque is Figure 1e.

the beta and bottom interaction terms. In the Northern Hemisphere, this cluster covers the southern extent of the subpolar gyres. A zonal streak crosses the equator in both the Atlantic and Pacific, but is absent in the Indian Ocean. The Southern Hemisphere has large Cluster 2 expanses in both the Pacific and Atlantic, but again not in the Indian Ocean.

The 15% of the ocean area selected by Cluster 3 are illustrated in Figure 3e (Quasi-Sverdrupian “negative wind torque/beta effect”). In Sverdrup balance, the wind torque and the beta effect are the only important terms. This theoretical relation is seen in Cluster 3 in the subtropical gyres where they are expected, together with regions in the Southern Ocean that the classical theory did not consider. Dominant areas in the sub-

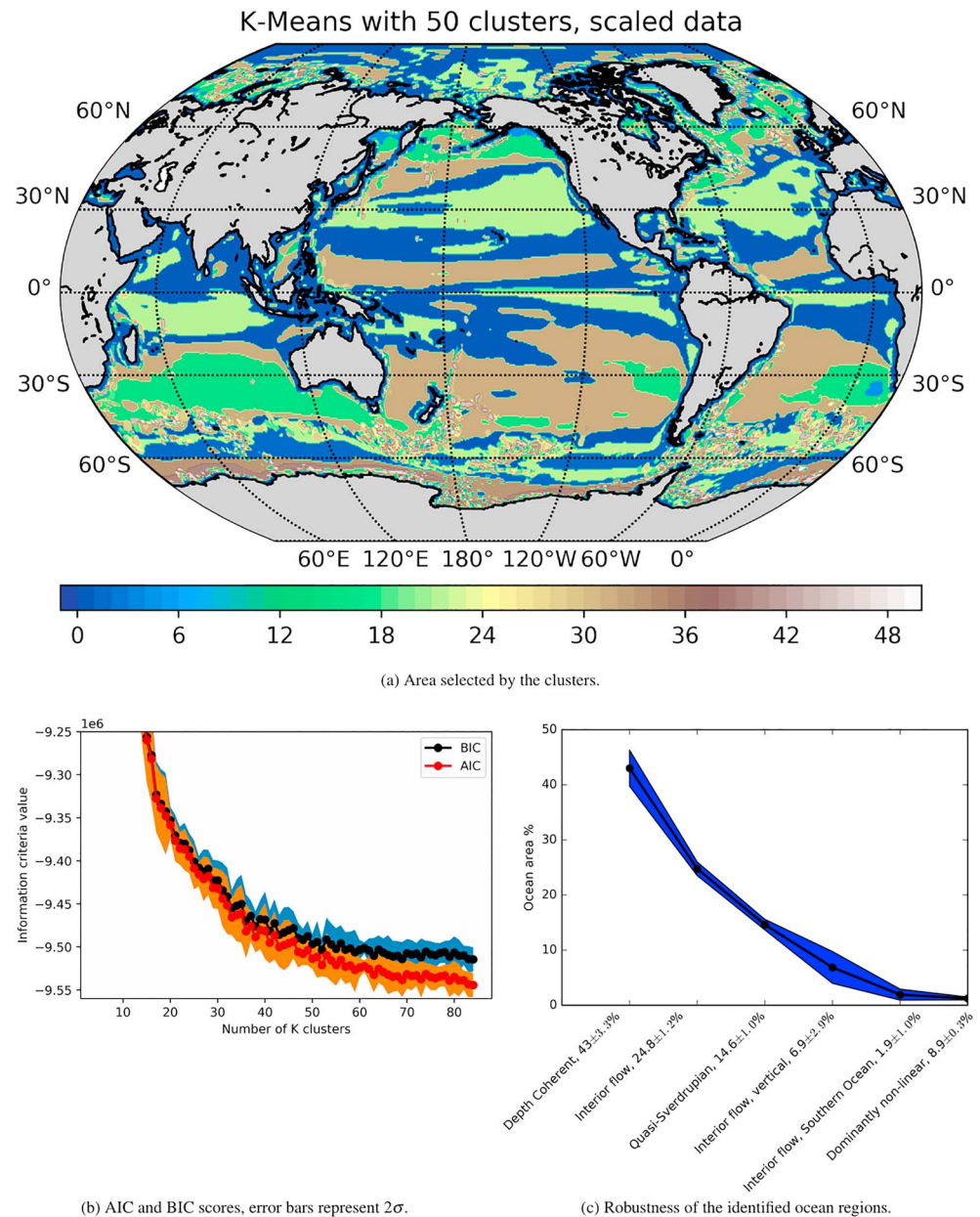


Figure 2. Top figure (a) illustrates the area selected by the clusters. Colors represent clusters in arbitrary order. The depth coherent ocean region (43%, Figure a) in dark blue, Interior flow (24.8%, Figure 3c) in light brown, Quasi-Sverdrupian (14.6%, Figure 3e) in light green, Interior flow, vertical, (6.9%, Figure 4a) in dark green, Interior flow, Southern Ocean, (1.9%, Figure 4c) in lighter blue, and the dominantly nonlinear torque cover remaining 8.9% (Figure 4e). Panel (b) illustrates that the Akaike Information Criteria (AIC) and Bayesian Information Criteria (BIC) asymptoting and a K of 50 is chosen for the analysis. Error bars of 2σ capturing the stochastic seed are shown. Panel (c) demonstrates the robustness of the algorithm with the ocean area, with 100 runs of the classification algorithm finding nearly identical areas (2σ error bars).

tropical gyres in the Northern Hemisphere Atlantic and Pacific stand out, together with thin streaks on the equator. Isolated streaks are seen in the Southern Ocean, and in a large area of the Southern Hemisphere tropical Indian Ocean. This region might be considered also as corresponding to quasi-Sverdrup balance.

Cluster 4 (Figure 4a: Interior flow, vertical, “positive wind torque, beta and bottom stress”), covers 7% of the ocean and is a complement to Cluster 1. In the Northern Hemisphere, the Cluster largely represents the northern edge of the subpolar gyre. In the Southern Hemisphere, it is found on the eastern edge of the Pacific and Atlantic basins, just to the south, and flaring out westward of the continental barrier. In the Indian

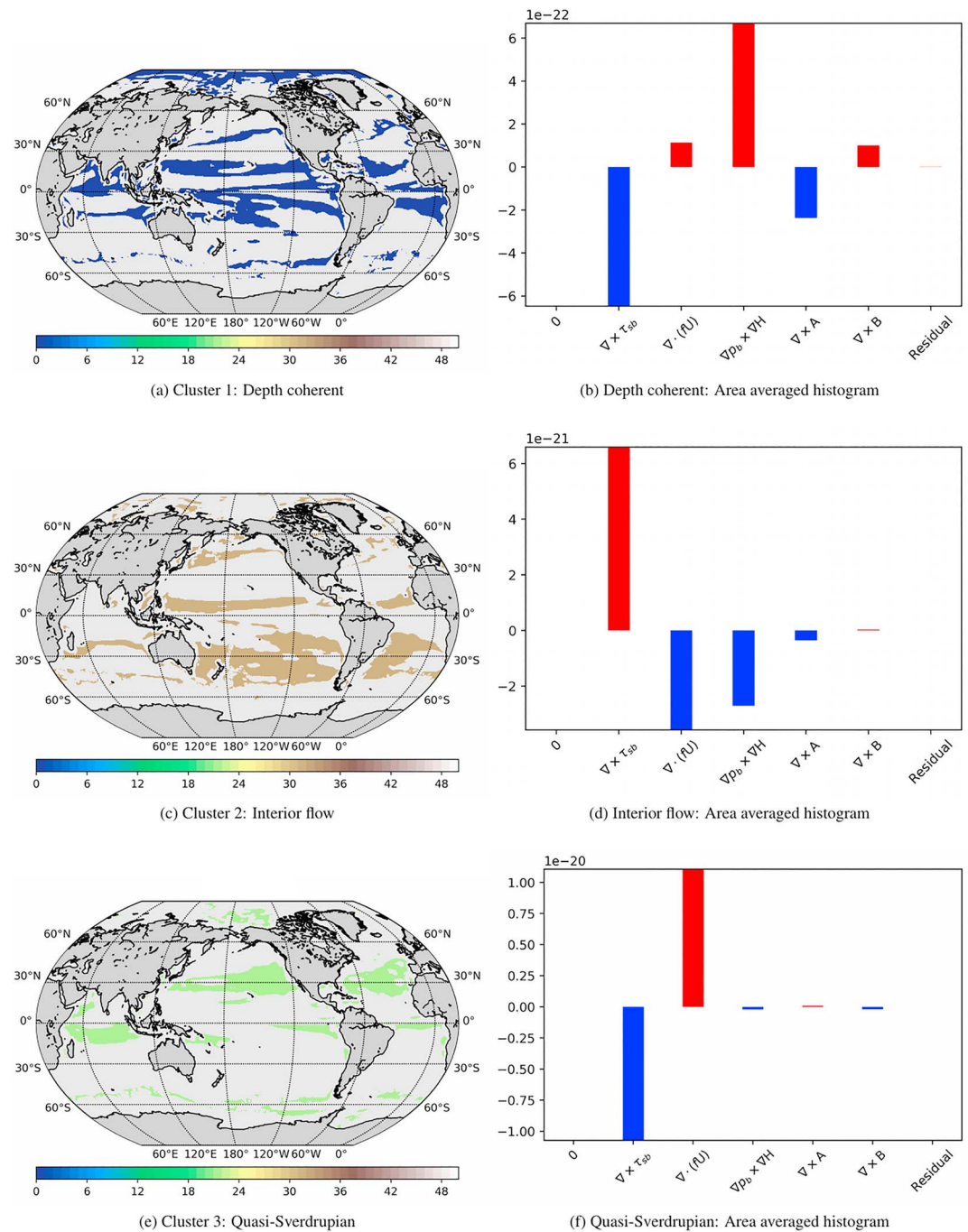


Figure 3. Maps of the selected locations (a, c, e) and corresponding area averaged histogram (b, d, f) of the terms in the barotropic vorticity equation. The color bar is kept, but the color/ordering of the map are arbitrary. Colors in the barchart indicate if barotropic vorticity is added (red) or removed (blue).

Ocean, this barrier can be seen to be New Zealand. The area of this dynamical regime fills the subtropical Indian Ocean down to the border with the Southern Ocean, where this regime is absent. Figure 4b illustrates that it is an amplified version of the dominant terms seen in Figure 3d, being an order of magnitude larger, but still having the wind as the major source of barotropic vorticity, with sinks in the Coriolis term and BPT. A small source exists in the nonlinear torque.

The “Southern Ocean gyre” is Cluster 5, covering only 2% of the global world ocean seen in Figure 4c. This Cluster is mainly seen in a series of zonal streaks in the Southern Ocean with negative wind torque and complements Cluster 4. Again, nonlinear torque is a small sink.

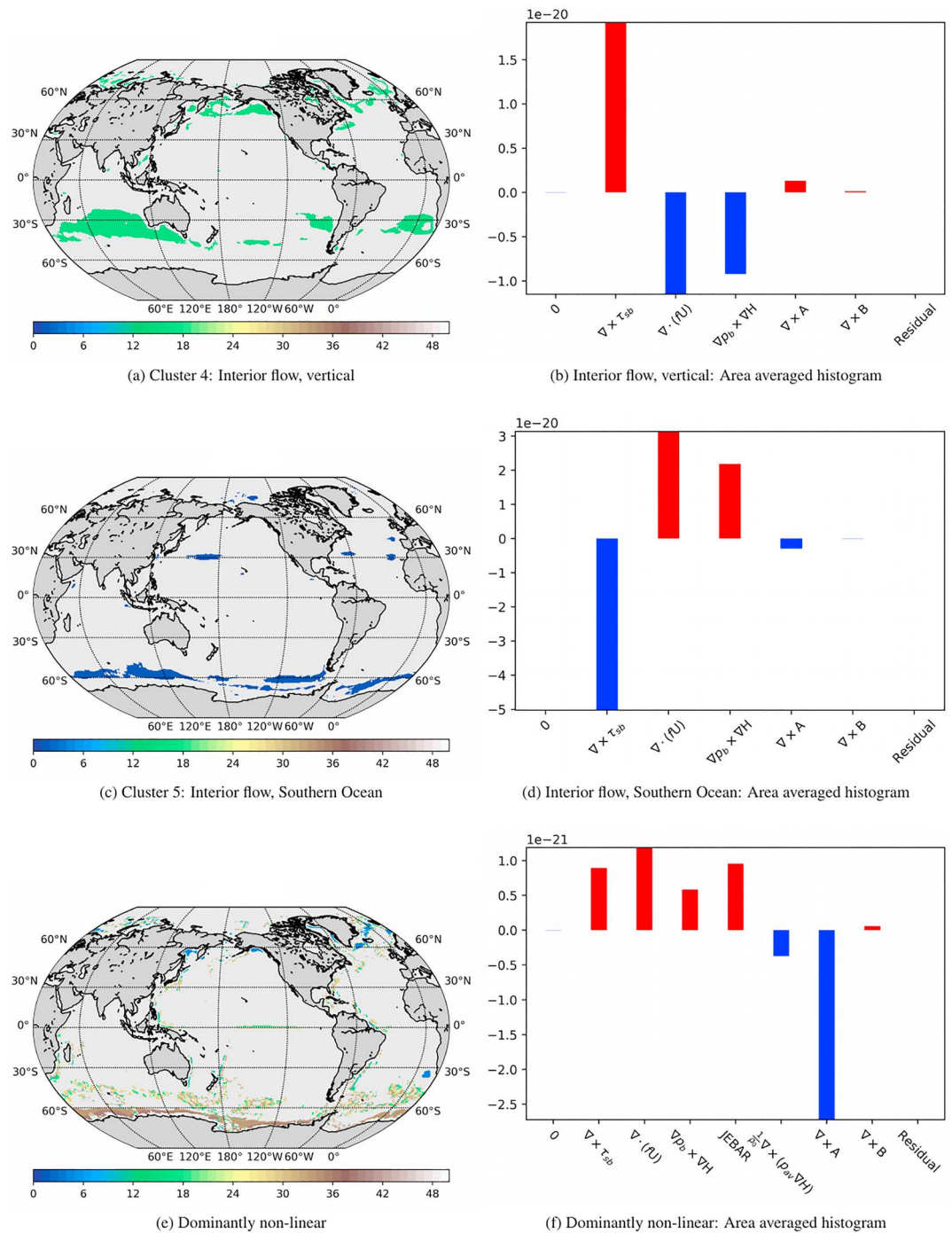


Figure 4. Maps of the selected locations (a, c, e) and corresponding area averaged histogram (b, d, f) of the terms in the barotropic vorticity equation. The colorbar is kept, but the color/ordering of the map are arbitrary. Colors in the barchart indicate if barotropic vorticity is added (red) or removed (blue).

A summary of the area of the remaining clusters that account for 9% of the world ocean is shown in Figure 4e: “Dominantly nonlinear.” Separate clusters have different colors. Areas of rough bathymetry stand out, such as the Pacific-Antarctic Ridge and the Drake Passage area. Figure 4f is the overall average, illustrating that the nonlinear contribution to the barotropic vorticity dominates, together with the Coriolis term. The different constituents are quite varied, but strong contributions from the nonlinear torque are consistently present. Their detailed discussion is the subject of a subsequent study.

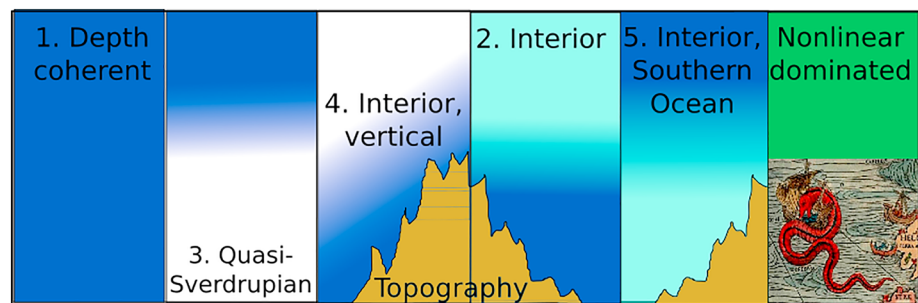


Figure 5. Schematic of identified regions with names and cluster numbers. The depth-coherent area implies a coherent vertical structure in Cluster 1. The quasi-Sverdrupian gyre in Cluster 3 is unique due to lack of bottom pressure torque (BPT). The interior flow, vertical, in Cluster 4 has a stronger momentum driven portion of the BPT, and topographic interactions begin to become important. The interior flow in Cluster 2 has a stronger baroclinic component to the BPT and feels topography. The interior flow, Southern Ocean, in Cluster 4 is like the interior flow in Cluster 2, but with contributions of opposite sign. The remainder is dominated by nonlinear contributions, and the barotropic interpretation is not appropriate.

4. Discussion and Conclusions

In the ECCOv4 state estimate, the barotropic vorticity equation closes very accurately. The 20-year time average ECCOv4 state estimate, was analyzed globally using *K*-means clustering to find regions dominated by groups of balanced terms in the BV equation. Large regions of the global ocean exhibit structures of consistent term balances as shown in Figure 5. Those balances vary among the wind stress, Coriolis, and BPT terms. Areas where the nonlinear torque are small suggest that the linearized BV is a good approximation. Areas where the nonlinear torque are important are found in western boundary regions, as well as the Southern Ocean where the Antarctic Circumpolar Current interacts with bathymetric obstacles. The momentum-dominated area (Cluster 1) implies a coherent vertical structure. Cluster 3 mainly in the subtropical gyre is unique in lacking significant contributions by BPT, implying it is shielded from topography. Transition zones have a stronger momentum-driven portion of the BPT and topographic interactions to become important. Cluster 4 has a stronger baroclinic component to the BPT, feeling topography. The Southern Ocean Cluster is like Cluster 2, but with contributions of opposite sign. Important nonlinear contributions are present in the remaining ocean, and the linearized barotropic interpretation is not appropriate.

In the North Atlantic Ocean, results are generally consistent with the inferences of the relative magnitude of the terms of the BV equation (Schoonover et al., 2016; Yeager, 2015). Cluster analysis reveals a shift from a barotropic flow in Clusters 1 and 3, to a strong interior flow (baroclinic meridional, North Atlantic Current, and North Atlantic Deep Water, and flow over the Mid Atlantic Ridge) in Clusters 2 and 4. Globally, the clustering illustrates that strong interior flow is present in vast expanses of the Southern Hemisphere, as well as in the North Pacific. Cluster 4 coincides with regions identified by Lumpkin and Speer (2007), Perez et al. (2011), and Speich et al. (2007) as areas of water mass transformation and intermediate pathways in the overturning circulation between surface and deep water. The quasi-Sverdrupian regime in Cluster 3 is not present in the South Atlantic and Pacific. In the Southern Ocean Clusters 3 and 5 dominate, with significant nonlinear contributions.

Five regions cover 91% of the world ocean. Residual areas collected here as “dominantly nonlinear” have a small spatial extent, but are dynamically important in the overall circulation. These regions include the Drake Passage region as well as the Antarctic-Pacific Ridge regions where the circulation interacts with topography and cross frontal transport likely takes place. In the Northern Hemisphere, areas in the Labrador Sea and on the continental shelf stand out as nonlinear. These nonlinear regions will be the subject of a separate study at higher resolution.

The sign and spatial distribution of the wind stress term suggests the importance of Ekman pumping (negative) or suction (positive). Equatorial and Southern Ocean regions show Ekman pumping, whereas the subpolar gyre areas have Ekman suction where mode waters are created. The BPT term mirrors the wind stress term, suggesting it acts as either a source or a sink in opposite complement to the wind stress. A lack of symmetry in wind driven gyres in the Southern and Northern Hemisphere show that the gyres are not

driven solely by the sign of the Ekman pumping. This complex relationship among the terms is the subject of future study.

The use of vertical integrals to describe the circulation is a simplification of what is a three-dimensional problem, as is the use of a model in which the important eddy field is only included through parameterizations such as in CMIP5 and CMIP6 (Church & Miles, 2013; Eyring et al., 2016; Stouffer et al., 2017). Future work is intended to apply this and related machinery to fully eddy-resolving ocean states.

Appendix A: K-Means and influence of Information Criteria

The K-means algorithm is related to methods such as Principal Component Analysis (PCA), more traditionally applied to oceanography. Where PCA attempts to represent all data vectors using a low-order combination of eigenvectors, minimizing the mean squared reconstruction error, the K-means algorithm represents the data vectors via a small number of clusters. This is also done to minimize the mean squared reconstruction error. In this manner, the K-means algorithm can be interpreted as a very sparse PCA.

Robustness of the regions in terms of the stochastic initialization is highlighted in Figure 2c, where the K-means clustering was run 100 times and mean and 2σ used as metrics in Table 1. The regimes identified are robust, with the extent of the subpolar gyre being the main area where the algorithm shows appreciable variance.

The K-means algorithm is initiated by scattering K first guesses of where the parameters/clusters could be. This initial guess introduces a stochastic element. The success of the algorithm is sensitive to K , as this determines how the hyperspace given by the dimensions is partitioned. As with regression analysis, adding parameters can increase the accuracy, but overfitting should be avoided. Determining the appropriate value of K , information criteria (AIC and BIC) are used to assess the quality of the statistical model. These measures weight the added accuracy with the cost of adding additional parameters, minimizing the expectation of the prediction error, are used:

$$\text{AIC} = 2K - 2 \ln(\mathcal{L}),$$

$$\text{BIC} = K \ln(n) - 2 \ln(\mathcal{L}),$$

where n is the number of data points and \mathcal{L} is the likelihood:

$$\mathcal{L} = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(\zeta_i - \hat{\zeta}_i)^2}{2\sigma^2} \right).$$

The parameter ζ_i is the observed, and $\hat{\zeta}_i$ is the prediction, so $(\zeta_i - \hat{\zeta}_i)^2$ are the prediction residuals. In the estimate, the AIC value is minimized, which determines the smallest appropriate order to represent the time series. As discussed by Priestley (1981) and Yang (2005), the AIC can overestimate the order. Figure 2b demonstrates that both the AIC and BIC stabilize at >35 K and the asymptotic nature of the regime.

The Euclidian distance is used, meaning the variance is assumed to be isotropic (meaning round). This leads to the standard practice of normalizing and standardizing data. To elucidate the impact of assumptions the algorithm makes for the classification, a more generalized form of clustering was also tested: Gaussian Mixture Models (GMM). GMM are used to assess the impact of assumptions relating to the covariance structure; spherical, diagonal, tied, or full covariance. Using the BIC, the results building on the BV data were not seen to be sensitive to this. However, this could be important at higher resolution as the K-means clustering problem is NP-hard and GMM could perform better.

References

- Adcroft, A., Hill, C., Campin, J. M., Marshall, J., & Heimbach, P. (2004). Overview of the formulation and numerics of the MIT GCM. In *Presented at the ECMWF Conference Proceedings* (pp. 139–150) Shinfield Park, Reading, UK.
- Akaike, H. (1973). "Information theory and an extension of the maximum likelihood principle". In B. N. Petrov & F. Csaki (Eds.), *Second International Symposium on Information Theory* (pp. 267–281). Tsahkadsor, Armenia, USSR, September 2–8, 1971, Budapest: Akademiai Kiad.
- Ardyna, M., Claustre, H., Sallee, J., D'Ovidio, F., Gentili, B., van Dijken, G., et al. (2017). Delineating environmental control of phytoplankton biomass and phenology in the Southern Ocean. *Geophysical Research Letters*, *44*, 5016–5024. <https://doi.org/10.1002/2016GL072428>

Acknowledgments

This work was funded by the U.S. National Aeronautics and Space Administration Sea Level Change Team (contract NNX14AJ51G) and through the ECCO Consortium funding via the Jet Propulsion Laboratory. M. S. acknowledges Anne Reinarz, Roosa Tikkanen, and Katherine Rosenfeld as well as the python scikit-learn toolbox. ECCOV4, release 2, model output is available in this website (ftp://mit.ecco-group.org/ecco_for_las/version_4/release2/). ECCOV4, release 2, model output is available in this website (ftp://mit.ecco-group.org/ecco_for_las/version_4/release2/). M. S. developed the concept, designed the method, and performed the numerical simulations. M. S. wrote the paper under the guidance of C. W. C. W. and P. H. contributed equally to the final version of the manuscript.

- Bishop, C. M. (2006). *Pattern recognition and machine learning (Information science and statistics)*. Berlin, Heidelberg: Springer-Verlag.
- Breuhl, S., Lofland, K. R., Semenchuk, E. M., Rokicki, L. A., & Penzien, D. B. (1999). Use of clusters analysis to validate IHS diagnostic criteria for migraine and tension-type headache. *Headache*, 39(3), 181–9. A study of validating diagnostic criteria using k-means on symptom patterns.
- Church, J. A., & Miles, E. R. (2013). Sea Level Change. In T. F. Stocker, D. Qin, G.-K. Plattner, M. Tignor, S. K. Allen, J. Boschung, et al. (Eds.), *Climate change 2013: The physical science basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* (Chap. 13, pp. 1137–1216). Cambridge: Cambridge University Press.
- ECCO Consortium (2017a). A twenty-year dynamical oceanic climatology: 1994–2013. Part 1: Active scalar fields: Temperature, salinity. Dynamic Topography, Mixed-Layer Depth, Bottom Pressure.
- ECCO Consortium (2017b). A twenty-year dynamical oceanic climatology: 1994–2013. Part 2: Velocities. Property Transports, Meteorological Variables, Mixing Coefficients.
- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E. (2016). Overview of the coupled model intercomparison project Phase 6 (CMIP6) experimental design and organization. *Geoscientific Model Development*, 9, 1937–1958. <https://doi.org/10.5194/gmd-9-1937-2016>
- Forget, G., Campin, J.-M., Heimbach, P., Hill, C. N., Ponte, R. M., & Wunsch, C. (2015). ECCO version 4: An integrated framework for non-linear inverse modeling and global ocean state estimation. *Geoscientific Model Development*, 8, 3071–3104.
- Fukumori, I., Heimbach, P., Ponte, R. M., & Wunsch, C. (2018). A dynamically-consistent ocean climatology. *Bulletin of the American Meteorological Society*, 99, 2107–2128. <https://doi.org/10.1175/BAMS-D-17-0213.1>
- Hauser, J., & Rybakowski, J. (1997). Three clusters of male alcoholics. *Drug Alcohol Depend*, 48(3), 243–250. An example of clustering behavior types in addiction research.
- Hughes, C. W., & Williams, S. D. P. (2010). The color of sea level: Importance of spatial variations in spectral shape for assessing the significance of trends. *Journal of Geophysical Research*, 115, C10048. <https://doi.org/10.1029/2010JC006102>
- Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. (2013). “Using ensemble of neural networks to learn stochastic convection parameterizations for climate and numerical weather prediction models from data simulated by a cloud resolving model”. *Advances in Artificial Neural Systems*, 2013, 13. <https://doi.org/10.1155/2013/485913>
- Kubat, M. (Ed.) (2015). *Introduction to machine learning* (2nd ed.). New York: Springer Publishers.
- Kulis, B., & Jordan, M. I. (2012). Revisiting k-means: new algorithms via Bayesian nonparametrics. In *Proceedings of the 29th International Conference on Machine Learning* (pp. 1–14). (ICML '12) July 2012 Edinburgh, UK 5135202-s2.0-84867132578.
- Lumpkin, R., & Speer, K. (2007). Global ocean meridional overturning. *Journal of Physical Oceanography*, 37, 2550–256.
- MacQueen, J. B. (1967). Some methods for classification and analysis of multivariate observations. In *Proceedings of 5-th Berkeley Symposium on Mathematical Statistics and Probability* (Vol. 1, pp. 281–297). Berkeley: University of California Press.
- Munk, W. (1950). On the wind-driven ocean circulation. *Journal of Atmospheric Sciences*, 7, 80–93. [https://doi.org/10.1175/1520-0469\(1950\)007<0080:OTWDOC>2.0.CO;2](https://doi.org/10.1175/1520-0469(1950)007<0080:OTWDOC>2.0.CO;2)
- Munk, W. H., & Palmén, E. (1940). Note on the dynamics of the Antarctic Circumpolar Current. *Tellus*, 3, 53–55.
- Myers, P. G., Fanning, A. F., & Weaver, A. J. (1996). JEBAR, bottom pressure torque, and gulf stream separation. *Journal of Physical Oceanography*, 26, 671–683.
- Perez, R., Garzoli, S., Meinen, C., & Matano, R. (2011). Geostrophic velocity measurement techniques for the meridional overturning circulation and meridional heat transport in the South Atlantic. *Journal of Atmospheric and Oceanic Technology*, 28, 1504–1521. <https://doi.org/10.1175/JTECH-D-11-00058.1>
- Priestley, M. B. (1981). *Spectral analysis and time series*. London: Academic Press.
- Reynolds, R. W., Chelton, D. B., Roberts-Jones, J., Martin, M. J., Menemenlis, D., & Merchant, C. J. (2013). Objective determination of feature resolution in two sea surface temperature analyses. *Journal of Climate*, 26, 2514–2533.
- Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling 2.0: A blueprint for models that learn from observations and targeted high-resolution simulations.
- Schoonover, J., Dewar, W., Wienders, N., Gula, J., McWilliams, J. C., Molemaker, M. J., et al. (2016). North Atlantic barotropic vorticity balances in numerical models. *Journal of Physical Oceanography*, 46(1), 289–303.
- Sonnevald, M., Wunsch, C., & Heimbach, P. (2018). Linear predictability: A sea surface height case study. *Journal of Climate*, 31, 2599–2611. <https://doi.org/10.1175/JCLI-D-17-0142.1>
- Speich, S., Blanke, B., & Cai, W. (2007). Atlantic meridional overturning circulation and the Southern Hemisphere supergyre. *Geophysical Research Letters*, 34, L23614. <https://doi.org/10.1029/2007GL031583>
- Stammer, D., Balmaseda, M., Heimbach, P., Koehl, A., & Weaver, A. (2016). Ocean data assimilation in support of climate applications: Status and perspectives. *Annual Review of Marine Science*, 8, 491–518.
- Stommel, H. (1948). The westward intensification of wind-driven ocean currents. *Eos, Transactions American Geophysical Union*, 29(2), 202–206. <https://doi.org/10.1029/TR029i002p00202>
- Stouffer, R. J., Eyring, V., Meehl, G. A., Bony, S., Senior, C., Stevens, B., & Taylor, K. E. (2017). CMIP5 scientific gaps and recommendations for CMIP6. *Bulletin of the American Meteorological Society*, 98, 95–105. <https://doi.org/10.1175/BAMS-D-15-00013.1>
- Wunsch, C., & Heimbach, P. (2007). Practical global oceanic state estimation. *Physica D*, 230, 197–208.
- Wunsch, C., & Heimbach, P. (2013). Dynamically and kinematically consistent global ocean circulation and ice state estimates, (2nd ed.). In G. Siedler, S. M. Griffies, J. Gould, & J. A. Church (Eds.), *Ocean Circulation and Climate* (Vol. 103, pp. 553–579). Academic Press.
- Xu, Y., & Fu, L. (2012). The effects of altimeter instrument noise on the estimation of the wavenumber spectrum of sea surface height. *Journal of Physical Oceanography*, 42, 2229–2233. <https://doi.org/10.1175/JPO-D-12-0106.1>
- Yang, Y. (2005). “Can the strengths of AIC and BIC be shared?” *Biometrika*, 92, 937–950.
- Yeager, S. (2015). Topographic coupling of the Atlantic overturning and gyre circulations. *Journal of Physical Oceanography*, 45, 1258–1284.