

¹ **Unsupervised learning classifies global ocean dynamical regions**

² Maike Sonnewald*, Carl Wunsch

³ *Massachusetts Institute of Technology, 77 Massachusetts Ave, Cambridge, MA 02139, USA,*

⁴ *Harvard University, 26 Oxford Street, Cambridge, MA 02138, USA and The University of Texas
5 at Austin, 201 East 24th Street, Austin, TX 78712, USA*

⁶ & Patrick Heimbach

⁷ *The University of Texas at Austin, 201 East 24th Street, Austin, TX 78712, USA and*

⁸ *Massachusetts Institute of Technology, 77 Massachusetts Ave, Cambridge, MA 02139, USA*

⁹ *Corresponding author address: Maike Sonnewald, Massachusetts Institute of Technology, 77
¹⁰ Massachusetts Ave, Cambridge, MA 02139, USA.

¹¹ E-mail: maike_s@mit.edu

ABSTRACT

12 Global ocean dynamically consistent regions are identified using a
13 barotropic vorticity framework for a twenty-year mean from the Estimating
14 the Circulation and Climate of the Ocean (ECCO) state estimate. Closure
15 of the barotropic vorticity budget provides both classification and dynamical
16 interpretation, highlighting where linear theory is appropriate and where
17 non-linear terms are significant. The unsupervised learning algorithm,
18 K-Means, demonstrates five unambiguous and basin independent global
19 regimes. A “Momentum dominated” area covers $57 \pm 2\%$ of the ocean area.
20 Surface and bottom stress terms are balanced there by the bottom pressure
21 torque and the non-linear terms. The “Transition zone” covers $18 \pm 0.7\%$,
22 and is found mainly in highly baroclinic subpolar regions. The surface
23 and bottom stress are largely balanced by the Coriolis term and the bottom
24 pressure torque. “Subtropical gyres” cover $11 \pm 0.5\%$, characterized by
25 a “Quasi-Sverdrupian” regime where the Coriolis term is balanced by the
26 wind and bottom stress term. The “Subpolar gyre” region characterised
27 by baroclinic dynamics covers $5.0 \pm 1.9\%$, and is a muted version of the
28 “Transition zone”. Similar gyre dynamics are seen in the “Southern Ocean
29 gyre” covering $1.4 \pm 0.7\%$. The shift from the subpolar through the transition
30 zone to the subtropical gyres is a shift from strong interior flow to one
31 dominated by e.g. the western boundary currents. Remaining $6 \pm 0.2\%$ of
32 the area is “Dominantly non-linear”, and found in key areas in the Southern
33 Ocean and along western boundaries.

35 **1. Motivation**

36 The global ocean contains diverse dynamic and kinematic regions. This diversity renders
37 difficult forming generalizations and understanding of the global system. A central purpose
38 of this paper is to objectively identify regions of common dynamics as a step towards their
39 understanding. For purposes of exploring appropriate methodologies, the study is restricted to the
40 barotropic vorticity balance of a time-mean global circulation, although it is readily generalized
41 to many other ocean circulation characteristics.

42

43 Characterizing the global ocean and reducing the overwhelming dimensionality of observed
44 motions has a long history in oceanography. Using theoretical conjectures, Stommel (1948) and
45 Munk (1950) assumed topography is not important in the subtropics, while topography has been
46 key in understanding Southern Ocean dynamics (Munk and Palmén 1951). From theoretical
47 conjectures and numerical models, it is not obvious what level of complication is merited, or if
48 any useful global description is plausible. General Circulation Models (GCMs) solve the primitive
49 equations in detail, and using the success of simple characterizations of ocean features can be
50 assessed against them. The need for a full GCM approach can be tested regionally, and its success
51 and failure understood in terms of captured key dynamics.

52

53 To this end, an “unsupervised learning classification” algorithm is applied here to a depth-
54 integrated, barotropic vorticity (BV) equation. Adequacy of the BV framework is assessed for
55 decadal timescales. The success of work such as Munk (1950) relies on a depth-integrated
56 view, where lateral friction provides closure in what is often referred to as Munk boundary layer
57 structure for Western Boundary Currents (WBC). A quasi-Sverdrupian framework was invoked

58 by him, where the depth-integrated meridional flow is given by the local wind-stress curl. The
59 potential success of such frameworks suggests predictive skill based purely on the wind field,
60 and that a barotropic, depth-integrated view is appropriate to characterize the subtropical ocean.
61 GCM studies and observational work present a more intricate ocean, where baroclinic structure
62 and non-linearities are key to ocean dynamics. If this complexity is correct for large regions of
63 the ocean, or if they govern key dynamics on relevant timescales, the simplified theoretical view
64 would be inadequate.

65

66 The presence of dynamical regimes is suggested by the structures of the wind-stress forcing
67 and the geometry of the ocean basins including the underlying topography. Classifying and
68 identifying regions in the world ocean is done here by using the BV equation and showing that the
69 global budget can be closed. Yeager (2015); Schoonover *et al.* (2016); Le Bras *et al.* (in review)
70 have assessed the dynamics of the BV budget focusing on the North Atlantic. The link to the
71 overturning and gyre circulation, as well as to the importance of mechanisms such as buoyancy
72 forcing were discussed. Here the procedures are global.

73

74 Much previous work assessing global ocean dynamics has been constrained to the surface.
75 Xu and Fu (2012) used surface quasigeostrophic theory and sea surface height (SSH) data from
76 altimeters to show differing regimes of geostrophic turbulence. Major current systems have
77 been identified, but with altimeters alone, the application to a comprehensive assessment of
78 global dynamical regimes is limited. Hughes and Williams (2010); Sonnewald *et al.* (2018) use
79 linear statistical models, finding global patterns in SSH. Those global patterns are connected to
80 circulation patterns and planetary waves, and the statistical approach using autoregressive and
81 moving average models was suggestive of global regimes. The present work extends the pattern

82 determination methodology.

83

84 Objective methods to identify patterns in data are common in many fields ranging from
85 pharmaceutical to engineering applications (Kulis and Jordan 2012; Breuhl *et al.* 1999; Hauser
86 and Rybakowski 1997). Applications in the sciences have been investigated both in the prognostic
87 and diagnostic sense (Krasnopolsky *et al.* 2013; Schneider *et al.* 2017).

88

89 Objective classification from machine learning has the potential to powerfully simplify the
90 interpretation of ocean dynamics. The objective aspect allows an unbiased assessment of the
91 success of the theoretical framework used, and the dimensionality reduction offered by using a
92 classifier is a potentially very useful approach in physical oceanography. Ardyna *et al.* (2017)
93 applied a similar method to identify regions with distinct biological activity, and Liang *et al.*
94 (in review) used K-Means clustering to identify key regions for data collection to build maps
95 of nitrate in the Southern Ocean. Compared to the more familiar principal component/factor
96 analysis (PCA), the K-Means cluster analysis has an additional “categorical” constraint. PCA and
97 K-Means maximize the same objective function, but are not identical as discussed in the appendix.

98

99 **2. Methods**

100 *a. ECCOv4 State Estimate*

101 The BV equation is used to assess the extent to which we can close this budget in the version
102 4 Estimating the Circulation and Climate of the Ocean (ECCOv4) release 2 estimate described
103 by Wunsch and Heimbach (2013); Forget *et al.* (2015) and others (see also ECCO Consortium

104 (2017a,b)). The state estimate is global 1° with tropical and high latitude mesh refinement. A
 105 least-squares with Lagrange multipliers approach is used to obtain the state estimate. The result is
 106 a *free-running* version of the MIT General Circulation Model (MITgcm, Adcroft *et al.* (2004))—
 107 with adjusted input variables—solving the primitive equations. In contrast to most “reanalysis”
 108 products, the ECCO oceanic state satisfies basic conservation laws for enthalpy, salt, volume, and
 109 momentum remaining largely within error estimates of a diverse set of global data (Wunsch and
 110 Heimbach 2007, 2013; Stammer *et al.* 2016). Regions without data are filled in a dynamically con-
 111 sistent way using the dynamics, avoiding the use of untested statistical hypotheses e.g., Reynolds
 112 *et al.* (2013).

113 *b. Barotropic Vorticity*

114 To classify the ocean’s dynamical regions the BV equation is used, one based on the solutions
 115 to the primitive equations from the ECCOv4 state estimate. Illustrating the simplification, the
 116 familiar momentum equations of an ocean on a thin shell on a rotating sphere are:

$$\partial_t \mathbf{u} + f \mathbf{k} \times \mathbf{u} = -\frac{1}{\rho_0} \nabla p + \frac{1}{\rho_0} \partial_z \tau + \mathbf{a} + \mathbf{b}, \partial_z p = -g\rho, \nabla \cdot \mathbf{v} = 0. \quad (1)$$

117 The pressure, acceleration due to gravity, density and vertical shear stress are denoted p , g , ρ
 118 and τ respectively, with ρ_0 the reference density; the three dimensional velocity field $\mathbf{v} = (\mathbf{u},$
 119 $\mathbf{v}, \mathbf{w}) = (\mathbf{u}, \mathbf{w})$; the unit vector is denoted \mathbf{k} ; planetary vorticity as a function of latitude ϕ in
 120 $f\mathbf{k} = (0, 0, 2\Omega \sin \phi)$; the viscous forcing by vertical shear is denoted $\partial_z \tau$; the non-linear terms
 121 are \mathbf{a} and the horizontal viscous forcing \mathbf{b} includes subgrid-scale parameterizations. Assuming a
 122 steady state, the vertical integral from the surface $z = \eta(x, y, t)$ to the water depth below the surface
 123 $z = H(x, y)$ gives

$$\beta V = \frac{1}{\rho_0} \nabla p_b \times \nabla H + \frac{1}{\rho_0} \nabla \times \tau + \nabla \times \mathbf{A} + \nabla \times \mathbf{B}, \quad (2)$$

124 where $\nabla \cdot \mathbf{U} = 0$, $\mathbf{U} \cdot \nabla f = \beta V$, the bottom pressure is denoted p_b , $\mathbf{A} = \int_H^\eta \mathbf{ad}z$ and $\mathbf{B} = \int_H^\eta \mathbf{bd}z$.
 125 The LHS of equation (2) is the planetary vorticity advective term, while the RHS of equation (2)
 126 is the bottom pressure torque (BPT), the wind and bottom stress curl, the non-linear torque and
 127 the viscous torque, respectively. The non-linear torque is composed of three terms:

$$\nabla \times \mathbf{A} = \nabla \times \left[\int_{-H}^0 \nabla \cdot (\mathbf{u}\mathbf{u}) dz \right] + [w\zeta]_{z=-H}^{z=0} + [\nabla w \times \mathbf{u}]_{z=-H}^{z=0}. \quad (3)$$

128 The RHS on equation (3) represents the curl of the vertically integrated momentum flux di-
 129 vergence, the non-linear contribution to vortex tube stretching and the conversion of vertical
 130 shear to barotropic vorticity. Horizontal viscous forcing includes that induced by subgrid-scale
 131 parameterizations. Twenty-year averaged fields of the BV equation are used after a Laplacian
 132 smoother is applied.

133

134 c. Unsupervised learning: K-Means clustering

135 Assessing the presence of dominant global patterns of dominant terms in the BV equation, K-
 136 Means clustering is used. In that approach, no groups are prescribed and the data are not con-
 137 strained, something that otherwise is a first step in oceanographic analysis. K-Means clustering
 138 originates in signal processing, and is used to determine if discernible patterns exist in the vari-
 139 ations of the data. Patterns of variability constitute areas in a function space, given by the terms
 140 of the BV equation, that are more densely populated than others. If groups of dominant terms are
 141 present that differ significantly, they are robustly identified and are called “clusters”. The spatially
 142 varying terms in the BV equation were scaled to have zero mean and unit variance globally. The

¹⁴³ algorithm involves an interative minimization of the sum of squares of the Euclidian distance given
¹⁴⁴ by the partitioning of the hyperspace:

$$J = \sum_{j=1}^K \sum_{i=1}^n \|x_i^j - c_j\|^2 \quad (4)$$

¹⁴⁵ where the number of K (c_j) is a parameter fixed a priori, randomly scattered among the n data
¹⁴⁶ points (x). The distance between a data point x_i^j and c_j is given by $\|x_i^j - c_j\|^2$. Initially, each data
¹⁴⁷ point is associated with the closest K cluster. The position of c_j is recalculated, and the association
¹⁴⁸ reassessed. This procedure is iteratively repeated. The problem is computationally NP-hard, and
¹⁴⁹ sensitive to the initialization and choice of K. As detailed in the Appendix, the appropriate value
¹⁵⁰ of K is determined as $K > 35$ using the Akaike and Bayesian Information Criteria (AIC and BIC).
¹⁵¹ The AIC and BIC indicate robust regimes as they both have asymptote, suggesting no information
¹⁵² is gained increasing K. A fixed K=50 is used for the remaining analysis. The physical meaning of
¹⁵³ the K parameter is associated to the variance structure of the data. The algorithm partitions the
¹⁵⁴ parameter subspace using linear hyperplanes. The linearity constraint means that higher numbers
¹⁵⁵ of K can both assist in partitioning the subspace more appropriately, and also isolate noise. The
¹⁵⁶ appendix demonstrates the small sensitivity to the initial random seed of the algorithm, as well as
¹⁵⁷ the impact of varying K.

¹⁵⁸

¹⁵⁹ 3. Results

¹⁶⁰ Figure 1 illustrates the closure in ECCOv4 for the 20-year average of the BV terms in equation
¹⁶¹ (2). Individual terms are of order $\pm 10^{-9} ms^{-1}$ and the residuals have magnitudes $<< 1\%$ for 36%
¹⁶² of the ocean—a very small value (Figure 1), less than $\pm 10^{-12} ms^{-1}$. These small residuals permit
¹⁶³ going forward with confidence. Some numerical issues do exist on the shelf and in shallow water

164 generally, but these only amount to 2.9% of the area of the global ocean and will be ignored.

165

166 Figure 1b illustrates where the Coriolis term is important in ECCOv4 for the 20-year average
167 from equation (2). This term is balanced by the wind and bottom stress BV terms shown in
168 Figure 1d (the bottom stress term is small) and BPT shown in Figure 1c. The remainder is largely
169 found in the non-linear BV contributions seen in Figure 1e, with the lateral viscous dissipation
170 largely being an order of magnitude smaller, apart from in localized regions in the Southern Ocean.

171

172 The wind and bottom stress BV term in Figure 1d are largely zonally symmetric in the 20-year
173 average, with large patterns of negative BV to the south in the Southern Ocean, large gyre patterns
174 visible in the Pacific and Atlantic basins. The BPT term in Figure 1c is associated with interactions
175 with steep bathymetry. For example, in the Southern Ocean a large positive patch leads towards
176 the Antarctic-Pacific ridge, and with a negative patch beyond. This structure is consistent with
177 vortex stretching as the ridge is crossed, and a similar feature is associated with the Mid-Atlantic
178 Ridge in the Atlantic sector of the Southern Ocean. Along Western Boundaries, BPT is positive to
179 the west and negative just adjacent to the east, consistent with studies such as Myers *et al.* (1996).
180 The BV of the non-linear terms is similarly concentrated along the western edge of basins where
181 WBCs are found, but it is less spatially coherent than the BPT term. The Southern Ocean stands
182 out as a region of high activity, particularly in the Atlantic sector. Lateral viscous dissipation in
183 the BV equation is comparatively small, but also concentrated in the Southern Ocean.

184

185 Figure 2 illustrates the spatial extent of the dynamical regimes picked out by the K-Means
186 algorithm. The numbering on the colorbar is arbitrary, and the structure is mainly found in
187 five geographically coherent patches named in Table 1. We use these names in the remaining

188 description. A zonal pattern suggests the wind-stress contribution to the BV is key, but much
189 zonal spatial structure appears in the Southern Ocean and along the western boundaries. Figures
190 3 and 4 isolate the geographical area and the associated area-averaged magnitudes of the terms
191 in the BV equation, ordered in decreasing percentage coverage of the global ocean. The bar
192 charts represent the area-averaged balance of the terms in the BV equation from the spatial areas
193 determined as distinct by the clustering. Table 1 summarizes the results, listing the corresponding
194 figure and dominant terms in the BV equation found in clusters.

195

196 Figure 3a illustrates the region dominated by a balance of the surface and bottom stress
197 terms covering 56% of the global ocean (Cluster 1). This balance suggests a barotropic regime
198 dominated by an apparent “Depth coherent” structure from surface to bottom, and it is found in
199 zonal streaks in the tropics, and a thin ribbon in the Southern Ocean notably in the Pacific sector.
200 In the Northern Hemisphere, Cluster 1 areas surround the subtropical and subpolar gyres. Large
201 areas of the Arctic Seas are also in this Cluster. Figure 3b demonstrates the balance of terms that
202 made the K-Means algorithm pick out the various regions. Wind and bottom stress terms are the
203 major source of negative barotropic vorticity, while the BPT adds positive vorticity. Non-linear
204 terms add negative vorticity.

205

206 Figure 3c illustrates the spatial region covering the next largest dynamical region covering 18%
207 of the ocean area (Cluster 2: “Transition zone”). In the Northern Hemisphere, this cluster covers
208 the southern region of the subpolar gyres. A zonal streak crosses the equator in both the Atlantic
209 and Pacific, but is absent in the Indian Ocean. The Southern Hemisphere has large Cluster 2
210 expanses in both the Pacific and Atlantic, but again not in the Indian Ocean. The bar chart in
211 Figure 3d highlights, as expected, the wind as the major source of barotropic vorticity, with sinks

212 in the Coriolis term and BPT. A small sink appears in the non-linear terms.

213

214 Figure 3e illustrates the 11% of the ocean area selected by the next Cluster (3: “Subtropical
215 gyre”). The subtropical gyres in the Northern Hemisphere Atlantic and Pacific stand out, together
216 with thin streaks on the Equator. Isolated streaks are seen in the Southern Ocean, and in a large
217 area of the Southern Hemisphere tropical Indian Ocean. Figure 3f shows a dominant balance
218 between the input of barotropic vorticity from the Coriolis term and sinks in the wind and bottom
219 stress. This balance corresponds to a quasi-Sverdrupian regime in the depth integrated ocean.

220

221 Figure 4a shows the area covered by Cluster 4 covering 7% of the ocean (“Subpolar gyre”).
222 This Cluster is largely a complimentary poleward extension to the cluster covering 18% of the
223 ocean seen in Figure 3c. In the Northern Hemisphere, the Cluster largely represents the northern
224 edge of the subpolar gyre. In the Southern Hemisphere, it is found on the eastern edge of the
225 Pacific and Atlantic basins, just to the south, and flaring out westwards of the continental barrier.
226 In the Indian Ocean, this barrier can be seen to be New Zealand or Australia, and the area of this
227 dynamical regime fills the subtropical Indian Ocean down to the border with the Southern Ocean,
228 where this regime is absent. Figure 4b illustrates that it is an amplified version of the dominant
229 terms seen in figure 3d, being an order of magnitude larger, but still having the wind as the major
230 source of barotropic vorticity with sinks in the Coriolis term and BPT. A small source exists in
231 the non-linear terms.

232

233 The Southern Ocean is better represented in the area covering 2% of the global world ocean
234 seen in Figure 4c (Cluster 5: “Southern Ocean gyre”), as seen mainly in a series of streaks in the
235 Southern Ocean, near 60°S where there is no continental block. Isolated areas are also seen in the

236 Northern Hemisphere. Figure 4d illustrates that this region, Cluster 5 is approximately an inverse
237 of the barcharts for the area representing the Northern extension of the subtropical gyres covering
238 7% of the world ocean. The wind is the major sink of barotropic vorticity, with sources in the
239 Coriolis term and BPT. Again, non-linear terms are a small sink.

240

241 Figure 4e is a summary of the area of the remaining terms that account for 6.3% of the
242 world ocean (“Dominantly non-linear”). Areas of rough bathymetry stand out, such as the
243 Pacific-Antarctic Ridge and the Drake Passage area. Figure 4f illustrates that the non-linear
244 contribution to the barotropic vorticity dominates, together with the Coriolis term. The different
245 constituents are quite varied, but strong contributions from the non-linear terms are consistently
246 present. Such regions will be discussed in a subsequent paper.

247

248 4. Discussion and Conclusions

249 The barotropic vorticity (BV) equation in a 20-year average state estimate is analyzed for
250 the world ocean. Figure 1 shows that the global ocean has large regions of coherent dynamical
251 term balances as displayed in Figure 5. Those balances are nuanced among the wind-stress,
252 Coriolis and bottom pressure torque (BPT) terms. Areas where the non-linear terms are small
253 suggest that the linearized BV is a good approximation. Areas where the non-linear terms are
254 important are found in western boundary regions, as well as the Southern Ocean where the
255 Antarctic Circumpolar Current interacts with bathymetric obstacles. The momentum dominated
256 area implies a coherent vertical structure. The subtropical gyre is unique in lacking significant
257 contributions by BPT. The transition zone has a stronger momentum driven portion of the BPT,
258 and topographic interactions begin to become important. The subpolar gyre has a stronger

259 baroclinic component to the BPT and feels topography. The Southern Ocean gyre is like the
260 subpolar one, but with contributions of opposite sign. The remaining ocean has important by
261 non-linear contributions, and the linearized barotropic interpretation is not appropriate.

262

263 Consistent with the Yeager (2015) analysis in the North Atlantic Ocean, there is a region with
264 a shift from strong interior flow (baroclinic meridional, North Atlantic Current, NAC, and North
265 Atlantic Deep Water, NADW, flow over the Mid Atlantic Ridge) to a more barotropic flow (deep
266 western boundary current, DWBC). The geographical extent and transition between the 18%
267 (Transition zone) and 11% (subtropical) areas make sense as they are mirror images of each
268 other, but with the difference in wind sign and BPT interactions. The 5% (Subpolar gyre) region
269 complements the subtropical region, being an intensified extension having small contributions
270 of non-linear terms. In the global picture, strong interior flow is present in vast expanses of
271 the Southern Hemisphere, as well as in the North Pacific. The quasi-Sverdrupian regime in
272 the Northern Hemisphere subtropics is virtually absent in the South Pacific, consistent with the
273 absence of a strong barotropic flow structure.

274

275 The relative sizes of the Clusters need not imply their relative importance in the global
276 circulation. Areas collected here as “Dominantly non-linear” have a small spatial extent, but are
277 found in regions that are known to be important for ocean dynamics. The Drake Passage region as
278 well as the Antarctic-Pacific Ridge are areas where the circulation interacts with topography, and
279 these regions are places where non-linear effects are large. Similarly in the Northern Hemisphere,
280 areas in the Labrador sea and shelf stand out as non-linear.

281

282 The sign and spatial distribution of the wind-stress term suggests the importance of Ekman
283 pumping(negative)/suction(positive). The equatorial and Southern Ocean regions show Ekman
284 pumping, while the subpolar gyre areas have Ekman suction where mode waters are created. The
285 BPT term mirrors the wind-stress term, suggesting it acts as either a source or a sink in opposite
286 complement to the wind-stress. A full description of the dynamics is outside the scope of the
287 current study.

288

289 K-Means clustering is not the only method for determining regions of dynamical commonality;
290 generalised gaussian mixture models are another possibility. Also, other covariance shapes can
291 be explored, with some initial results in figure 8. Results are little changed, but they may be an
292 important element in moving to higher resolution estimates.

293

294 *Summary*

295 The barotropic vorticity equation closes very accurately in the 20-year time-average ECCOv4
296 state estimate. Balances observed involve the three dominant terms of the wind, Coriolis and
297 bottom pressure torque terms, with significant covariance and input from the non-linear terms
298 in some regions. Clusters of similar dynamics were identified using the K-Means clustering
299 method. Five regions cover 93% of the world ocean: depth coherent (57%), transition zone (18%),
300 subtropical gyre (11%), subpolar gyre (5%) and the Southern Ocean gyre (1.4%). The residual
301 area is dominated by the non-linear terms and will be the subject of a future study.

302

303 **5. Acknowledgments**

304 This work was funded by the US National Aeronautics and Space Administration Sea Level
305 Change Team (contract NNX14AJ51G) and through the ECCO Consortium funding via the Jet
306 Propulsion Laboratory. MS acknowledges the advice and support of Anne Reinarz, Edward Dod-
307 dridge, Roosa Tikkanen and Katherine Rosenfeld.

308 **K-Means and influence of Information Criteria**

309 The K-Means algorithm is related to methods such as PCA, more traditionally applied to
310 oceanography. Where PCA attempts to represent all data vectors using a low order combination
311 of eigenvectors, minimizing the mean squared reconstruction error, the K-Means algorithm
312 represents the datavectors via a small number of clusters. This is also done to minimize the mean
313 squared reconstruction error. In this manner, the K-Means algorithm can be interpreted as a very
314 sparse PCA.

315

316 The K-Means algorithm is initiated by scattering K first-guesses of where the parameters/clusters
317 could be. This initial guess introduces a stochastic element. The success of the algorithm is
318 sensitive to K , as this determines how the hyperspace given by the dimensions is partitioned. As
319 with regression analysis, adding parameters can increase the accuracy, but overfitting should be
320 avoided. To determine the appropriate value of K , the Akaike and Bayesian Information Criteria
321 (AIC and BIC) are used, minimizing the expectation of the prediction error:

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

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

322 where n is the number of datapoints and \mathcal{L} is the likelihood:

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

323 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
324 estimate, the AIC value is minimized, which determines the smallest appropriate order to represent
325 the time-series. As discussed by Priestley (1981) and Yang (2005), the AIC can overestimate the
326 order. Figure 6 demonstrates that both the AIC and BIC stabilise at $> 35K$, and the asymptotic
327 nature of the regime.

328

329 Robustness of the regions in terms of the stochastic initialisation is highlighted in Figure 7,
330 where the K-Means clustering was run 100 times. The mean and 2σ are used in Table 1. The
331 regimes identified are robust, with the extent of the subpolar gyre being the main area where the
332 algorithm shows appreciable variance.

333

334 To elucidate the impact of assumptions the algorithm makes for the classification, a more
335 generalised form of clustering is also tested: Gaussian Mixture Models. For context, K-Means
336 clustering is the limit where all covariances are diagonal, small and equal. Gaussian Mixture
337 Models are used to assess the impact of assumptions relating to the covariance structure; spherical,
338 diagonal, tied or full covariance are assessed as seen in Figure 8. Using the BIC to assess the
339 impact, this is found to likely to much change present results, but it could be important at higher
340 resolution high resolution, as the K-Means clustering problem is NP-hard.

341

342 References

- 343 Adcroft, A., Hill, C., Campin, J. M., Marshall, J., & Heimbach, P. (2004). Overview of the formulation and numerics of the MIT GCM (pp.
344 139-150). Presented at the ECMWF Conference Proceedings, Shinfield Park, Reading, UK.
- 345 Akaike, H. (1973), "Information theory and an extension of the maximum likelihood principle", in Petrov, B.N. ; Csiki, F., 2nd International
346 Symposium on Information Theory, Tsahkadsor, Armenia, USSR, September 2-8, 1971, Budapest: Akademiai Kiad, p. 267-281.
- 347 Ardyna, M., H. Claustre, J. Sallee, F. D'ovidio, B. Gentili, G. van Dijken, F. D'Ortenzio, and K. R. Arrigo, 2017: Delineating environmental control
348 of phytoplankton biomass and phenology in the southern ocean. *Geophys. Res. Lett.*, 44, 50165024, doi:10.1002/2016GL072428.
- 349 Breuhl S, et al. (1999). Use of clusters analysis to validate IHS diagnostic criteria for migraine and tension-type headache. *Headache*; 39(3):181-9.
- 350 A study of validating diagnostic criteria using k-means on symptom patterns.
- 351 Church, J.A., et al., 2013: Sea Level Change. In: Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the
352 Fifth Assessment Report of the Intergovernmental Panel on Climate Change [Stocker, T.F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J.
353 Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley (eds.)], Chapter 13, pp. 1137-1216, Cambridge University Press.
- 354 ECCO Consortium, 2017a, A Twenty-Year Dynamical Oceanic Climatology: 1994-2013. Part 1: Active Scalar Fields: Temperature, Salinity,
355 Dynamic Topography, Mixed-Layer Depth, Bottom Pressure.
- 356 ECCO Consortium, 2017b, A Twenty-Year Dynamical Oceanic Climatology: 1994-2013. Part 2: Velocities, Property Transports, Meteorological
357 Variables, Mixing Coefficients.
- 358 Forget, G., J.-M. Campin, P. Heimbach, C. N. Hill, R. M Ponte, and C. Wunsch, ECCO version 4: an integrated framework for non-linear inverse
359 modeling and global ocean state estimation, *Geo. Sci. Model Dev.*, 8, 2015
- 360 Hauser J and Rybakowski J (1997). Three clusters of male alcoholics. *Drug Alcohol Depend*; 48(3):243-50. An example of clustering behavior
361 types in addiction research.
- 362 Hirschi, J. J.-M., Blaker, A. T., Sinha, B., Coward, A., de Cuevas, B., Alderson, S., and Madec, G.: Chaotic variability of the meridional overturning
363 circulation on subannual to interannual timescales, *Ocean Sci.*, 9, 805-823, doi:10.5194/os-9-805-2013, 2013.
- 364 Chris W. Hughes, Simon D. P. Williams, The color of sea level: Importance of spatial variations in spectral shape for assessing the significance of
365 trends, *Journal of Geophysical Research*, 2010, 115, C10

366 Vladimir M. Krasnopolksy, Michael S. Fox-Rabinovitz, and Alexei A. Belochitski, "Using Ensemble of Neural Networks to Learn Stochastic
367 Convection Parameterizations for Climate and Numerical Weather Prediction Models from Data Simulated by a Cloud Resolving Model,"
368 Advances in Artificial Neural Systems, vol. 2013, Article ID 485913, 13 pages, 2013. doi:10.11552013485913

369 Kulis, B., Jordan, M. I.Revisiting k-means: new algorithms via Bayesian nonparametrics, Proceedings of the 29th International Conference on
370 Machine Learning (ICML '12)July 2012 Edinburgh, UK5135202-s2.0-84867132578

371 Le Bras, I., Toole, J.M., Sonnewald, M.: A bulk Potential Vorticity budget for the western North Atlantic based on observations. In review.

372 Liang, Y.-C. *et al.*, in review.

373 J. B. MacQueen (1967): "Some Methods for classification and Analysis of Multivariate Observations, Proceedings of 5-th Berkeley Symposium
374 on Mathematical Statistics and Probability", Berkeley, University of California Press, 1:281-297.

375 Munk, W., 1950: On the wind-driven ocean circulation. J. Meteor., 7, 7993, doi:10.1175/1520-0469(1950)007;0080:OTWDOC_2.0.CO;2.

376 Munk, W.H., and Palmén, E.; Note on the dynamics of the Antarctic Circumpolar Current, Tellus, 3, 53-55, 1940.

377 Myers, Paul G., Augustus F. Fanning, Andrew J. Weaver, 1996: JEBAR, Bottom Pressure Torque, and Gulf Stream Separation. J. Phys. Oceanogr.,
378 26, 671683. doi: [http://dx.doi.org/10.1175/1520-0485\(1996\)026<0671:JBPTAG>2.0.CO;2](http://dx.doi.org/10.1175/1520-0485(1996)026<0671:JBPTAG>2.0.CO;2)

379 Priestley, M. B. (1981) Spectral Analysis and Time Series, London: Academic Press.

380 Reynolds, R.W., D.B. Chelton, J. Roberts-Jones, M.J. Martin, D. Menemenlis, and C.J. Merchant, 2013: Objective Determination of Feature
381 Resolution in Two Sea Surface Temperature Analyses. J. Climate, 26, 2514-2533, <https://doi.org/10.1175/JCLI-D-12-00787.1>

382 Schoonover, J., Dewar, W., Wienders, N., Gula, J., McWilliams, J.C., Molemaker, M.J., Bates, S.C., Danabasoglu, G. and Yeager, S.: North
383 Atlantic Barotropic Vorticity Balances in Numerical Models, Journal of Physical Oceanography 2016 46:1, 289-303

384 Lan, S., Schneider, T., Stuartland, A., and Teixeira, J. (2017). Earth System Modeling 2.0: A Blueprint for Models That Learn From Observations
385 and Targeted High-Resolution Simulations.

386 Sonnewald, M., C. Wunsch, and P. Heimbach, 2018: Linear Predictability: A Sea Surface Height Case Study. J. Climate, 31, 25992611,
387 <https://doi.org/10.1175/JCLI-D-17-0142.1>

388 Stainforth, D., Aina, T., Christensen, C., Collins, M., Faull, N., Frame, D., Kettleborough, J., Knight, S., Martin, A., Murphy, J., Piani, C., Sexton,
389 D., Smith, L., Spicer, R., Thorpe, A. and Allen, M.: 2005, Uncertainty in predictions of the climate response to rising levels of greenhouse
390 gases, Nature 433(7024), 403-406. url:<http://oro.open.ac.uk/5025/>

- 391 D. Stammer, M. Balmaseda, P. Heimbach, A. Koehl, and A. Weaver, 2016: Ocean Data Assimilation in Support of Climate Applications: Status
392 and Perspectives. *Ann. Rev. Mar. Sci.*, 8, 491-518.
- 393 Stommel, H. (1948), The westward intensification of wind-driven ocean currents, *Eos Trans. AGU*, 29(2), 202206, doi:10.1029/TR029i002p00202.
- 394 Yang, Y. (2005), "Can the strengths of AIC and BIC be shared?", *Biometrika*, 92: 937-950, doi:10.1093/biomet/92.4.937.
- 395 Yeager, S.: Topographic coupling of the Atlantic overturning and gyre circulations, *J. Phys. Oceanogr.*, 45, 1258-1284. doi: <http://dx.doi.org/10.1175/JPO-D-14-0100.1>, 2015.
- 397 Wunsch, C., and P. Heimbach, 2007: Practical global oceanic state estimation. *Physica D*, 230, 197-208.
- 398 Wunsch, C. and P. Heimbach, 2013, Dynamically and kinematically consistent global ocean circulation and ice state estimates. In *Ocean Circulation*
399 and Climate 2nd Edition, Siedler et al., Eds.
- 400 Wunsch, C. The interpretation of short climate records, with comments on the North Atlantic and Southern Oscillations. *Bull. Am. Met. Soc.* 80,
401 245-255 (1999).
- 402 Wunsch, C. (2013), Covariances and linear predictability of the Atlantic Ocean, *Deep Sea Res., Part II*, 85, 228-243.
- 403 Wunsch, C. (2015), *Modern Observational Physical Oceanography: Understanding the Global Ocean*. Princeton University Press.
- 404 Xu, Y. and L. Fu, 2012: The Effects of Altimeter Instrument Noise on the Estimation of the Wavenumber Spectrum of Sea Surface Height. *J. Phys.*
405 *Oceanogr.*, 42, 22292233, <https://doi.org/10.1175/JPO-D-12-0106.1>

406 LIST OF TABLES

Table 1. Percentage of area covered by the area specific balance of the BV equation (??) and the corresponding map figure. Leading order terms are sorted by magnitude, colors indicating if barotropic vorticity is added (**red**) or removed (**blue**) 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.

21

Cluster	Area	Leading terms
1	57±1.8%, Depth coherent (Fig. 3a)	$\nabla \times \tau_{sb} + \nabla \times \mathbf{A} \approx \nabla p_b \times \nabla H$ (Fig. 3b)
2	18±0.7%, Transition zone (Fig. 3c)	$\nabla \times \tau_{sb} \approx \nabla p_b \times \nabla H + \nabla \cdot (\mathbf{fU})$ (Fig. 3d)
3	11±0.5%, Subtropical gyre (Fig. 3e)	$\nabla \times \tau_{sb} \approx \nabla \cdot (\mathbf{fU})$ (Fig. 3f)
4	5.0±1.9%, Subpolar gyre (Fig. 4a)	$\nabla \times \tau_{sb} \approx \nabla \cdot (\mathbf{fU}) + \nabla p_b \times \nabla H$ (Fig. 4b)
5	1.4±0.7%, Southern Ocean gyre (Fig. 4c)	$\nabla \times \tau_{sb} \approx \nabla \cdot (\mathbf{fU}) + \nabla p_b \times \nabla H$ (Fig. 4d)
6-50	5.7 ± 0.2%, Dominantly non-linear (Fig. 4e)	$\nabla \cdot (\mathbf{fU}) \approx \nabla \times \mathbf{A} + \nabla \times \tau_{sb}$ (Fig. 4f)

TABLE 1: Percentage of area covered by the area specific balance of the BV equation (2) and the corresponding map figure. Leading order terms are sorted by magnitude, colors indicating if barotropic vorticity is added (**red**) or removed (**blue**) 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.

413 LIST OF FIGURES

414	Fig. 1. The breakdown of the barotropic vorticity budget (ms^{-1}) over 1992-2013 in the 415 ECCOv4 State Estimate.	23
416	Fig. 2. The area selected by the clusters. The colors represent the clusters, and are in arbitrary order. 417 The momentum dominated ocean region (55.6%, Figure ??) is in dark blue, the Transition 418 zone (17.6%, Figure ??) is in light brown, the Subtropical Gyre (“Quasi-Sverdrup”, 11.2% 419 Figure?? in light green, the Subpolar gyre (6.6%, Figure ??) is in dark green, the Southern 420 Ocean gyre (2%, Figure ??) in lighter blue and the remaining areas where the non-linear 421 terms are large amounting to 6.3% of the world ocean and are presented in Figure ??.)	24
422	Fig. 3. Maps of the selected locations (left) and corresponding area averaged his- 423 togram (right) of the terms in the BV equation. The colorbar is kept, but the 424 color/ordering of the map are arbitrary. Colors in the barchart indicate if BV is 425 added (red) or removed (blue).	25
426	Fig. 4. Maps of the selected locations (left) and corresponding area averaged his- 427 togram (right) of the terms in the BV equation. The colorbar is kept, but the 428 color/ordering of the map are arbitrary. Colors in the barchart indicate if BV is 429 added (red) or removed (blue).	26
430	Fig. 5. Schematic of identified regions. The momentum dominated area implies a coherent verti- 431 cal structure. The subtropical gyre is unique due to lack of BPT. The transition zone has 432 a stronger momentum driven portion of the BPT, and topographic interactions begin to be- 433 come important. The subpolar gyre has a stronger baroclinic component to the BPT and 434 feels topography. The Southern Ocean gyre is like the subpolar, but with contributions of 435 opposite sign. The remainder is dominated by non-linear contributions, and the barotropic 436 interpretation is not appropriate.	27
437	Fig. 6. The AIC and BIC asymptote and we choose a K of 50 for our analysis. Error bars represent 438 2σ , capturing the stochastic start seen of the algorithm.	28
439	Fig. 7. The ocean area represented by the different areas from table ??, and 100 runs of the clas- 440 sification. Error bars represent 2σ , capturing the stochastic start seed of the classification 441 algorithm. The final point “Dominantly non-linear” represents the remainder in the 45 clus- 442 ters not shown.	29
443	Fig. 8. Trying different covariance models to check the convergence.	30

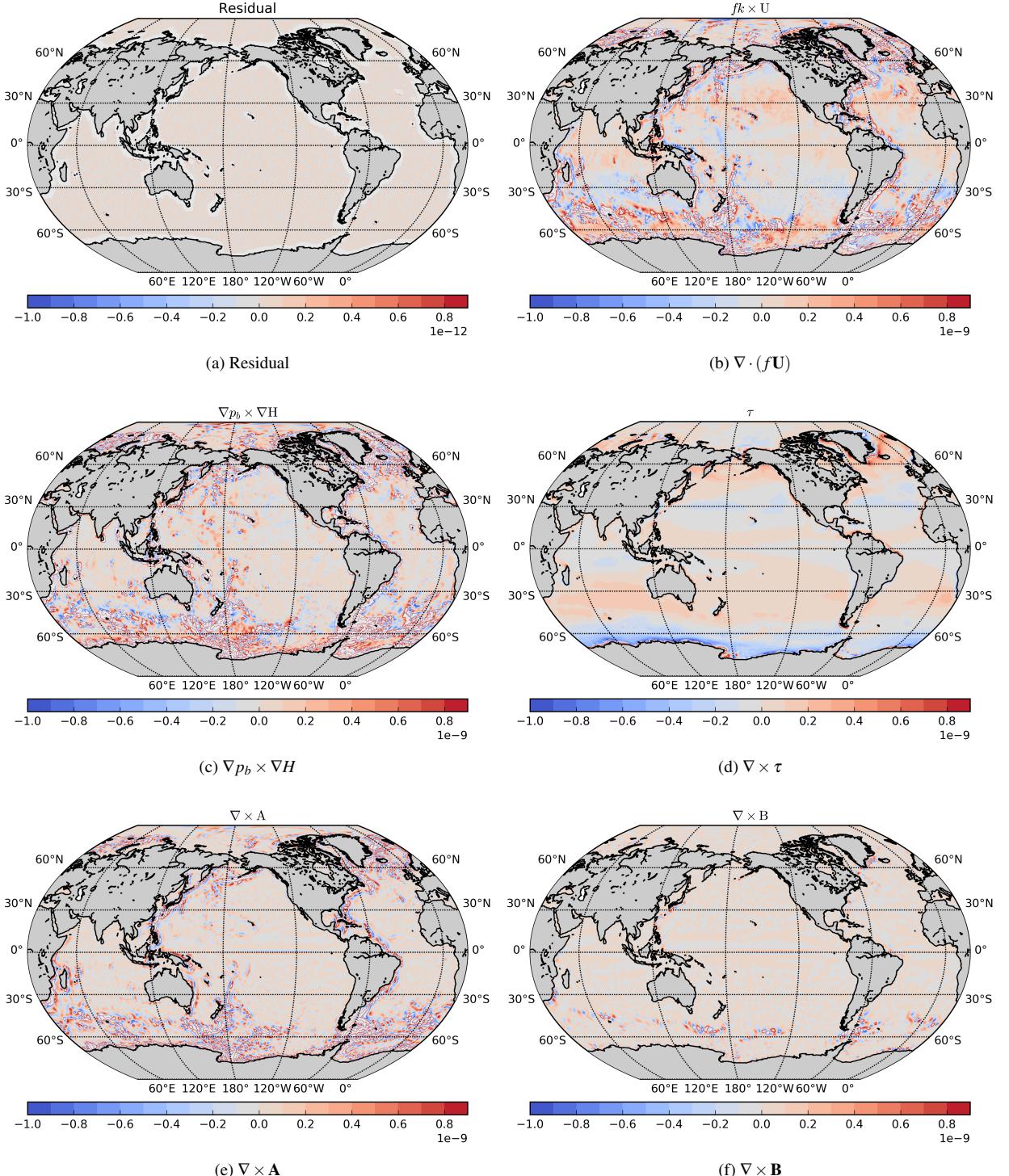


FIG. 1: The breakdown of the barotropic vorticity budget (ms^{-1}) over 1992-2013 in the ECCOv4 State Estimate.

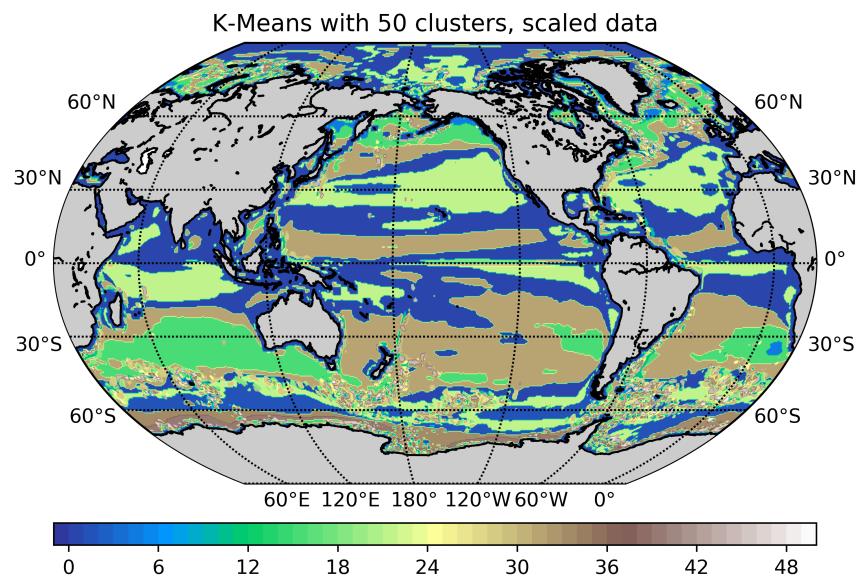
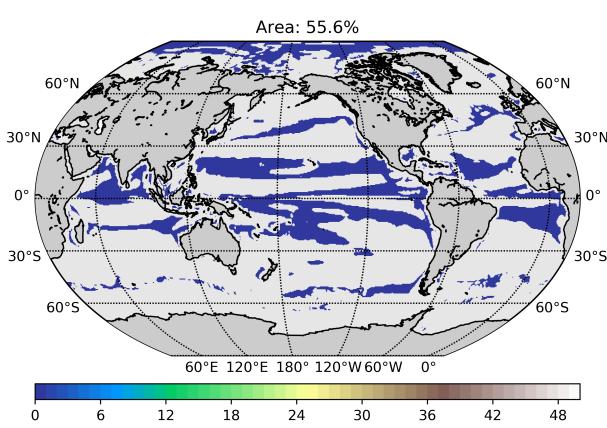
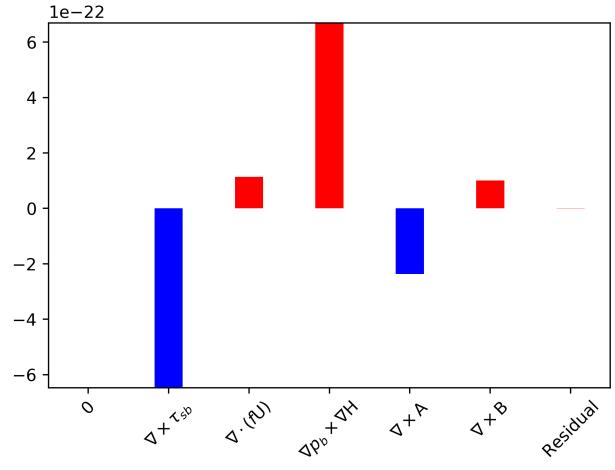


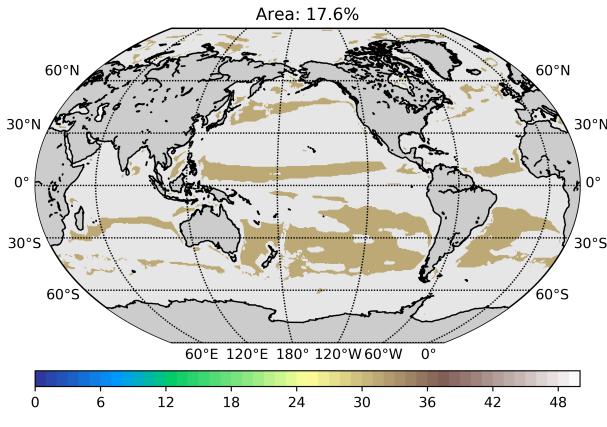
FIG. 2: The area selected by the clusters. The colors represent the clusters, and are in arbitrary order. The momentum dominated ocean region (55.6%, Figure 3a) is in dark blue, the Transition zone (17.6%, Figure 3c) is in light brown, the Subtropical Gyre (“Quasi-Sverdrup”, 11.2% Figure 3e in light green, the Subpolar gyre (6.6%, Figure 4a) is in gark green, the Southern Ocean gyre (2%, Figure 4c) in lighter blue and the remaining areas where the non-linear terms are large amounting to 6.3% of the world ocean and are presented in Figure 4e.)



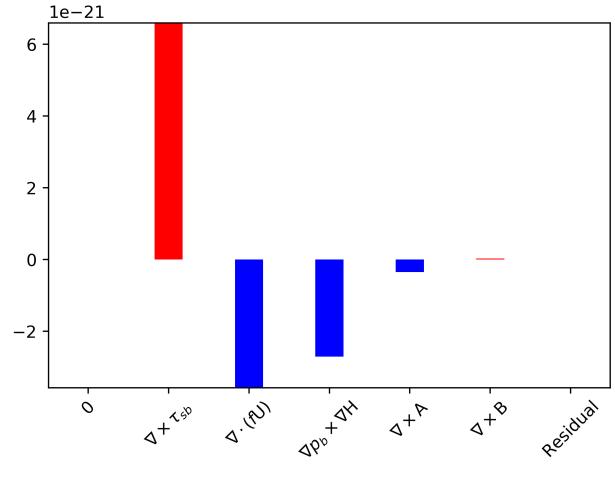
(a) Cluster 1: Depth coherent



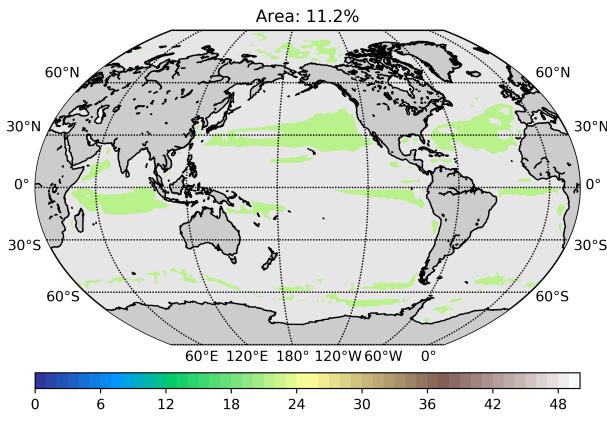
(b) Depth coherent: Area averaged histogram



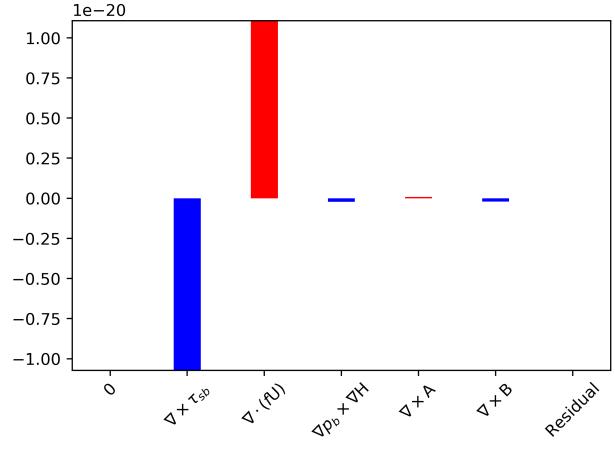
(c) Cluster 2: Transition Zone



(d) Transition zone: Area averaged histogram

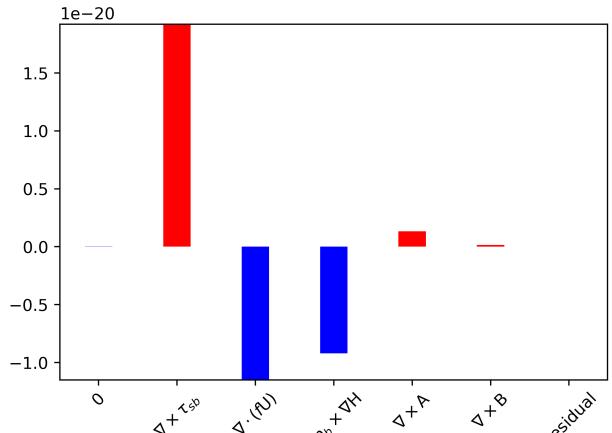
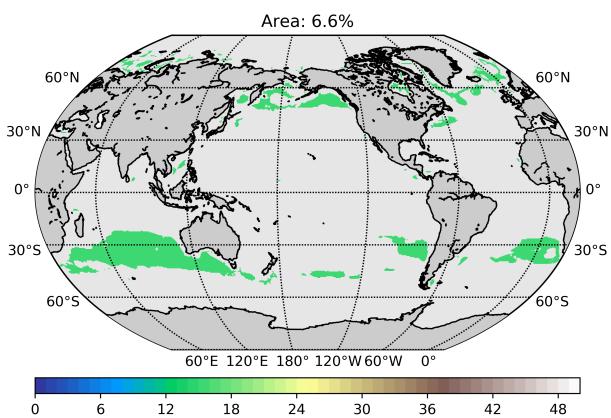


(e) Cluster 3: Subtropical gyre

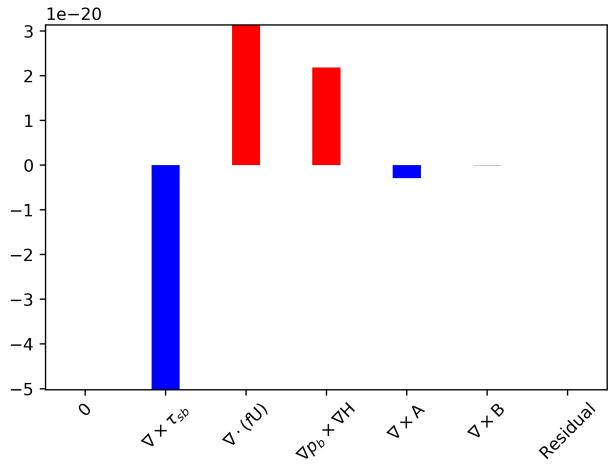
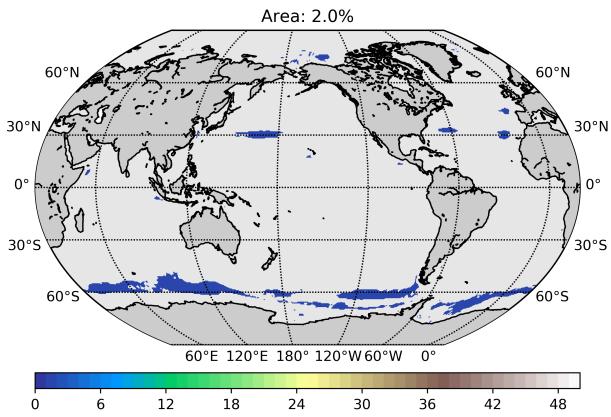


(f) Subtropical gyre: Area averaged histogram

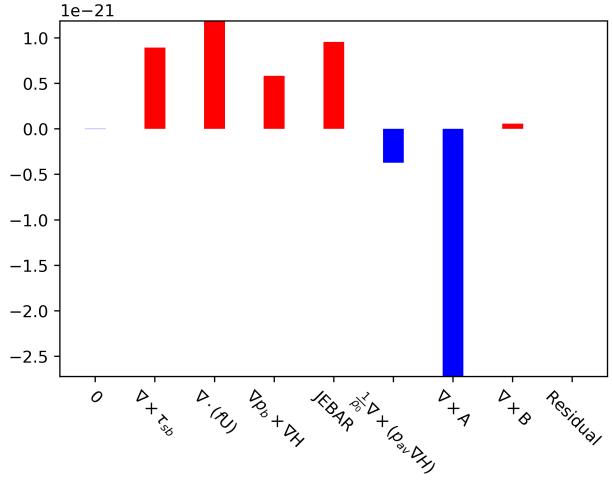
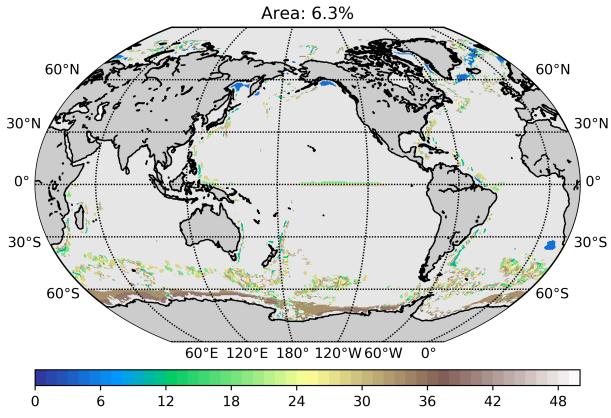
FIG. 3: Maps of the selected locations (left) and corresponding area averaged histogram (right) of the terms in the BV equation. The colorbar is kept, but the color/ordering of the map are arbitrary. Colors in the barchart indicate if BV is added (red) or removed (blue).



(b) Subpolar gyre: Area averaged histogram



(d) Southern Ocean gyre: Area averaged histogram



(f) Dominantly non-linear: Area averaged histogram

FIG. 4: Maps of the selected locations (left) and corresponding area averaged histogram (right) of the terms in the BV equation. The colorbar is kept, but the color/ordering of the map are arbitrary. Colors in the barchart indicate if BV is added (red) or removed (blue).

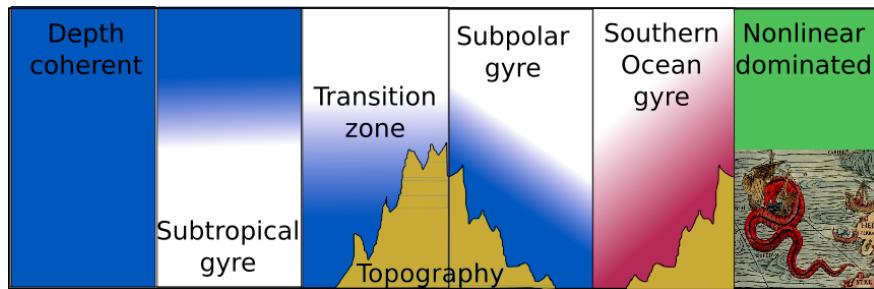


FIG. 5: Schematic of identified regions. The momentum dominated area implies a coherent vertical structure. The subtropical gyre is unique due to lack of BPT. The transition zone has a stronger momentum driven portion of the BPT, and topographic interactions begin to become important. The subpolar gyre has a stronger baroclinic component to the BPT and feels topography. The Southern Ocean gyre is like the subpolar, but with contributions of opposite sign. The remainder is dominated by non-linear contributions, and the barotropic interpretation is not appropriate.

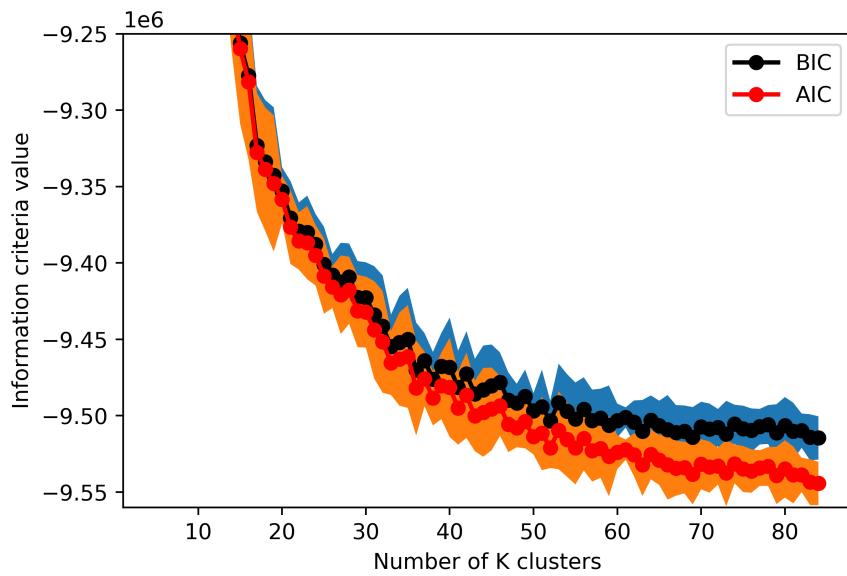


FIG. 6: The AIC and BIC asymptote and we choose a K of 50 for our analysis. Error bars represent 2σ , capturing the stochastic start seen of the algorithm.

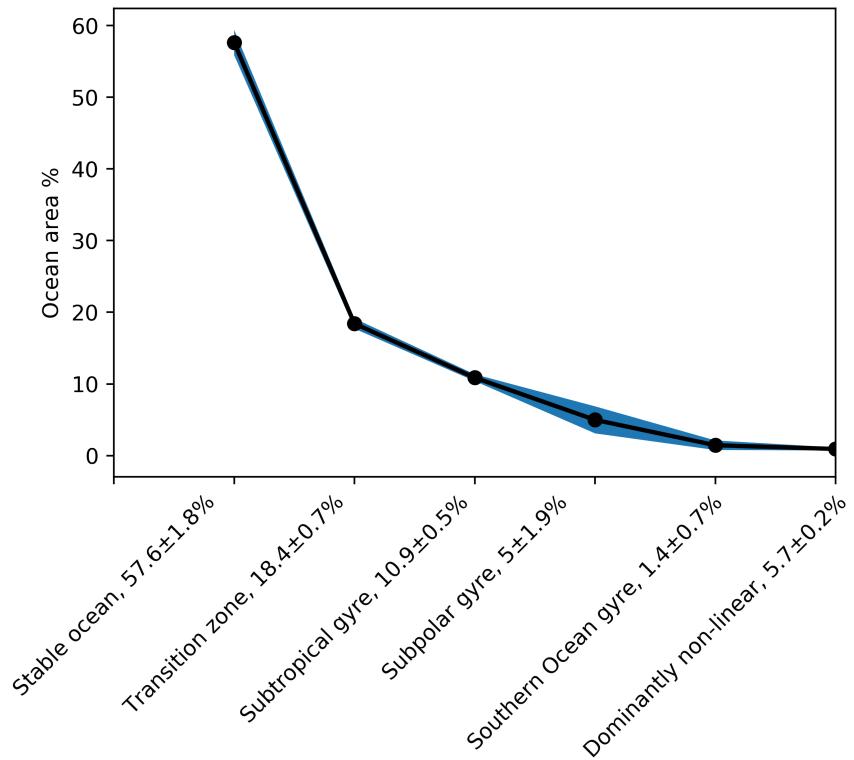


FIG. 7: The ocean area represented by the different areas from table 1, and 100 runs of the classification. Error bars represent 2σ , capturing the stochastic start seed of the classification algorithm. The final point “Dominantly non-linear” represents the remainder in the 45 clusters not shown.

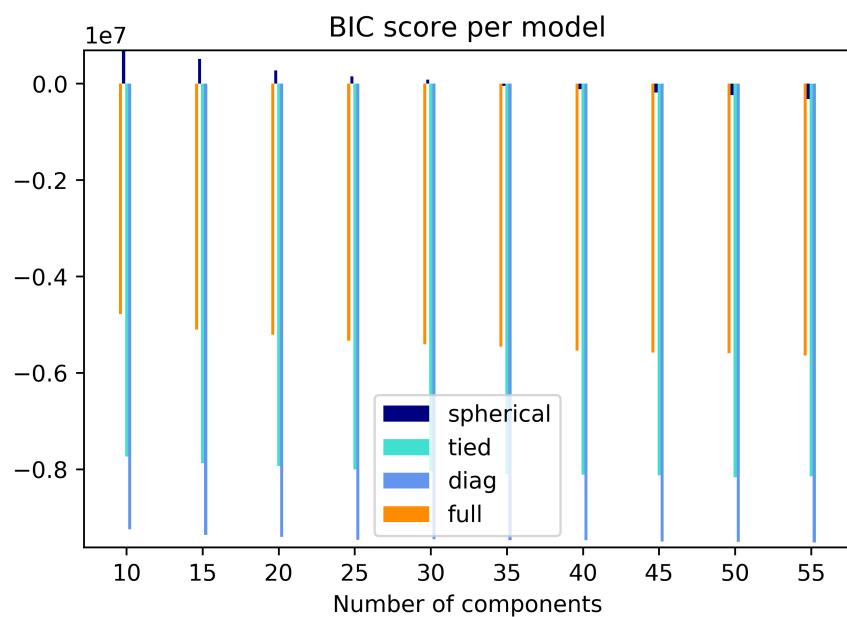


FIG. 8: Trying different covariance models to check the convergence.