

# Unsupervised Clustering of Southern Ocean Argo Float Temperature Profiles

Daniel C. Jones<sup>1</sup> , Harry J. Holt<sup>1,2</sup>, Andrew J. S. Meijers<sup>1</sup> , and Emily Shuckburgh<sup>1</sup> 

<sup>1</sup>British Antarctic Survey, Cambridge, UK, <sup>2</sup>Department of Physics, University of Cambridge, Cambridge, UK

## Key Points:

- We apply Gaussian mixture modeling (GMM) to Southern Ocean temperature data
- GMM identifies spatially coherent profile types without using latitude or longitude information
- GMM offers a complementary approach for objectively classifying temperature profiles

## Supporting Information:

- Supporting Information S1

## Correspondence to:

D. C. Jones,  
dannes@bas.ac.uk

## Citation:

Jones, D.C., Holt, H. J., Meijers, A. J. S., & Shuckburgh, E. (2019). Unsupervised clustering of Southern Ocean Argo float temperature profiles. *Journal of Geophysical Research: Oceans*, 124, 390–402. <https://doi.org/10.1029/2018JC014629>

Received 3 OCT 2018

Accepted 29 DEC 2018

Accepted article online 4 JAN 2019

Published online 17 JAN 2019

**Abstract** The Southern Ocean has complex spatial variability, characterized by sharp fronts, steeply tilted isopycnals, and deep seasonal mixed layers. Methods of defining Southern Ocean spatial structures traditionally rely on somewhat ad hoc combinations of physical, chemical, and dynamic properties. As a step toward an alternative approach for describing spatial variability in temperature, here we apply an unsupervised classification technique (i.e., Gaussian mixture modeling or GMM) to Southern Ocean Argo float temperature profiles. GMM, without using any latitude or longitude information, automatically identifies several spatially coherent circumpolar classes influenced by the Antarctic Circumpolar Current. In addition, GMM identifies classes that bear the imprint of mode/intermediate water formation and export, large-scale gyre circulation, and the Agulhas Current, among others. Because GMM is robust, standardized, and automated, it can potentially be used to identify structures (such as fronts) in both observational and model data sets, possibly making it a useful complement to existing classification techniques.

**Plain Language Summary** The Southern Ocean is an important part of the climate system, in part because it absorbs a large fraction of the heat and carbon that is added to the atmosphere/ocean system by human-driven fossil fuel burning. In this work, we use a machine learning technique to automatically sort Southern Ocean temperature measurements into groups based on how those temperature measurements change with depth. Different groups have the fingerprints of different large-scale circulation patterns, such as the powerful Antarctic Circumpolar Current that flows around Antarctica. The groups that we identify are consistent with our understanding of the Southern Ocean, which gives us confidence that our machine learning technique may be useful for automatically grouping measurements and computer model data in the future. This matters because the climate science community needs a new set of tools, possibly including the machine learning technique that we use in this paper, to deal with a very large, ever-increasing volume of observational and computer model data.

## 1. Introduction

The Southern Ocean (SO) is a critical component of Earth's climate system, having thus far absorbed greater than 75% of the energy added via anthropogenic emissions and 50% of the excess carbon (Fletcher et al., 2006; Frölicher et al., 2015). Its ability to absorb heat and carbon comes in part from its unique spatial structure and circulation, which features upwelling of cold, nutrient-rich waters and regions of dense water formation (Lumpkin & Speer, 2007). Characterizing and understanding the spatial variability of the SO remains an important and climatically relevant goal of modern oceanography.

Through decades of effort, the oceanographic community has converged on a description of ocean spatial variability that uses temperature, salinity, dynamical, and biogeochemical patterns to define different spatial structures (Emery, 2003; Talley, 2013, and references therein). For example, SO mode waters, located equatorward of the Antarctic Circumpolar Current (ACC), are commonly identified using potential vorticity minima and ranges of neutral density (Hanawa & Talley, 2001; Herraiz-Borreguero & Rintoul, 2011; Sallée et al., 2008, 2010a). Such systematic approaches employ the understanding that structural properties are “set” in their formation regions and modified by advection, mixing, and biogeochemical processes. In the SO, classification in latitude-longitude has traditionally been centered around several fronts of the ACC, defined by sharp transitions in sea surface height or neutral density (Kim & Orsi, 2014). The classical southern boundary (SBDY) of the ACC marks the transition between subpolar, gyre-dominated circulations and

lower latitude, more circumpolar flow. The ACC itself features three circumpolar fronts, namely, the Southern ACC Front (SACCF), the Polar Front (PF), and the Subantarctic Front (SAF; Orsi et al., 1995). These three fronts separate the subpolar SO from the subtropical domain (Naveira-Garabato et al., 2011).

The modern, property-driven classification scheme is extremely useful and will continue to be useful well into the future, but it is not necessarily ideal for every application. Many of the temperature, salinity, and density values used to delimit one structure from another are somewhat ad hoc and very specific (e.g., boundaries between different types of mode water). These schemes are useful for observational data analysis but difficult to apply to numerical models of the ocean, which do not necessarily feature exactly the same structure as the observed ocean (Sallée et al., 2013). In addition, traditional classification approaches that define structures by specific property ranges are limited by the fact that these properties may change over time, either as long-term trends (e.g., the warming of Antarctic Bottom Water [AABW] observed by Purkey & Johnson, 2010) or in terms of interannual variability (Naveira-Garabato et al., 2009). We suggest that it is prudent to develop and test alternative methods for the classification of oceanic temperature, salinity, and density structures, as a complement to existing expertise-driven methods.

Maze et al. (2017) have shown that Argo temperature profile data from the North Atlantic Ocean can be usefully grouped into classes using Gaussian mixture modeling (GMM), an unsupervised classification technique. GMM describes the spatial structure of Argo profiles as a collection of Gaussian modes whose means and standard deviations generally vary with pressure. In this work, we apply GMM to SO Argo temperature profiles in the upper 1,000 m of the water column. We find that GMM identifies several circumpolar classes, gyres, the Agulhas Current, and pathways broadly associated with the formation and export of mode and intermediate waters. In section 2, we describe the Argo data set and the basics of GMM. In section 3, we present the results of applying GMM to SO Argo data, and in sections 4 and 5 we offer a brief discussion and summarize our conclusions.

## 2. Methods

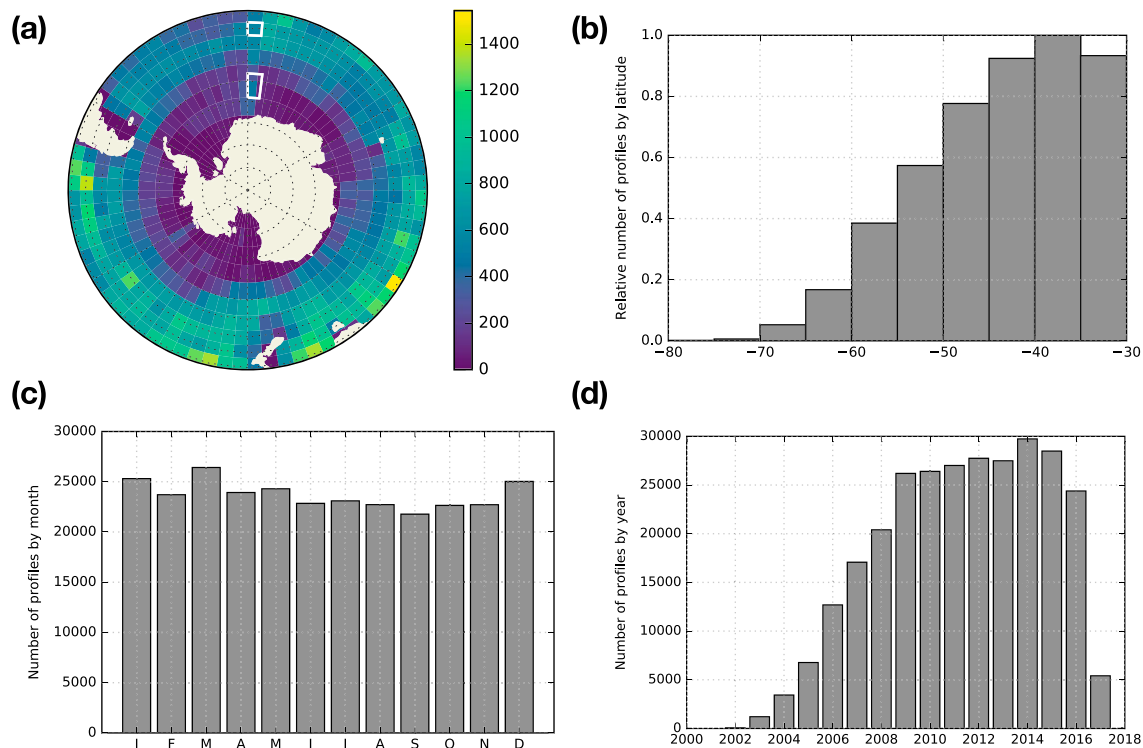
We applied an unsupervised classification method (i.e., GMM) to Southern Ocean Argo float data. In this section, we briefly describe the Argo data set and the basics of GMM. We use the *scikit-learn* machine learning library for Python (<http://scikit-learn.org/>), and the source code used for much of the analysis in this paper is available via Github (Holt & Jones, 2018). We refer the reader to Maze et al. (2017) for further details on applying GMM to Argo float data.

### 2.1. Argo Float Data Set

Argo floats are autonomous ocean instruments that measure, at minimum, the temperature and salinity of the ocean by periodically taking vertical profiles. Every 10 days, starting at a “neutral” position of 1,000 m, an Argo float dives down to 2,000 m before rising to the surface, taking a vertical profile of the water column along the way. The measurements are transmitted via satellite and are ultimately made freely available via the Argo Global Data Assembly Centers (GDACs) after some quality control checks. At the time of this writing, over 3,800 Argo floats are active in the global ocean, producing over 100,000 temperature and salinity profiles per year with an average spacing of 3° (<http://www.argo.ucsd.edu/>).

For this study, we selected all available Argo profiles south of 30°S that have been flagged by the GDACs as “observation good” (i.e., quality control flag = 1) covering the time period from 2001 to early 2017. More specifically, we used a vertically interpolated product with 400 equally spaced pressure levels ranging from 5 to 2,000 dbar in 5-dbar increments. After discarding profiles with greater than or equal to 6% NaN values (2% of the initial number of profiles) and discarding pressure levels with greater than or equal to 3% NaN values, we were left with 284,427 profiles, each with 192 pressure levels between 15 and 980 dbar. Most of these initially removed NaN values came from interpolation below roughly 1,000 dbar, as opposed to gaps in the original data set. We selected our NaN cutoff values based on the relatively large increase in the number of NaN values below 1,000 dbar. We replaced all remaining NaN values ( $\ll 1\%$  of the total temperature measurements) with linearly interpolated estimates using nearest neighbor values with respect to pressure. We refer to the resulting data set as the cleaned data set.

Because of the autonomous and free-drifting nature of the floats, the profiles are not distributed evenly in latitude/longitude (Figure 1). The profiles are more heavily concentrated in the Pacific sector (roughly 890 profiles per degree longitude, totalling 47% of profiles) and Indian sector (800 profiles per degree longitude,



**Figure 1.** Distribution of Argo temperature profiles from the cleaned data set. (a) Number of profiles in  $5^\circ \times 5^\circ$  bins. Two equal-area boxes are shown for reference (solid white lines). (b) Relative number of profiles by latitude, scaled by an area-weighting factor  $\cos(\varphi)$ , where  $\varphi$  is latitude. The temporal distribution of profiles shown by (c) month and (d) year.

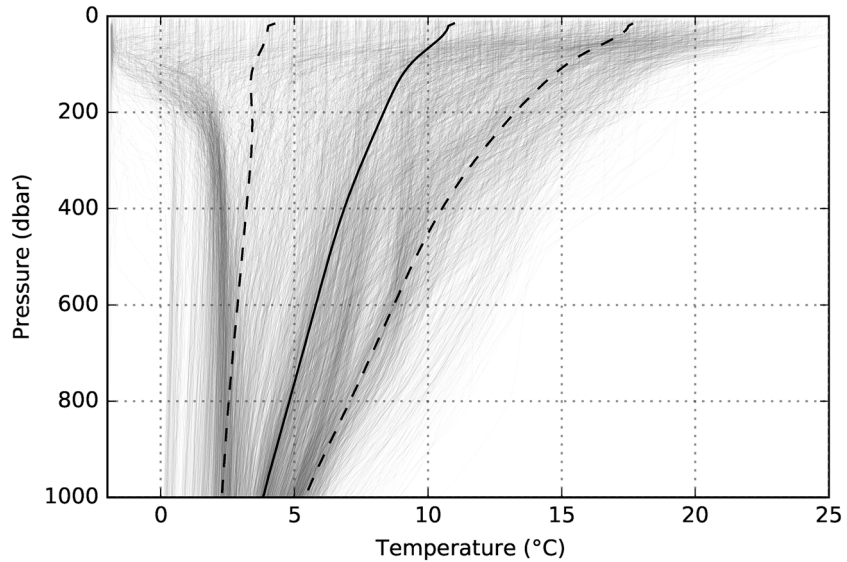
totalling 34% of profiles), with fewer profiles in the Atlantic sector (610 profiles per degree longitude, 19% of total). When counted in equal-area bins and plotted by latitude, we see that the number of profiles decreases toward Antarctica (Figure 1b), which is partly due to challenging operational conditions associated with seasonal sea ice, which can extend to just north of  $60^\circ\text{S}$  at maximum areal extent. The profiles are slightly overrepresented in the austral summer and autumn (December–February and March–May, 52% of profiles) and underrepresented in the Austral winter and spring (June–August and September–November, 48% of profiles), and the number of profiles increases until 2013 (Figures 1c and 1d). Since we selected an Argo data set that was created in early 2017, there are relatively few profiles from that year.

The profiles selected for this study display a large variety of vertical temperature structures (Figure 2). The range of temperatures is wider in the surface and considerably narrower with increasing pressure, in part reflecting the seasonal cycle in upper ocean temperatures. A large number of profiles feature colder temperatures near the surface and warmer temperatures in the interior, a physical arrangement that would be unstable to convection without the compensating effect of salinity. Profiles around Antarctica tend to be fresher at the surface and saltier in the interior due to glacial melt, freshwater flux, and the balance of evaporation/precipitation. This arrangement of temperature and salinity can be stable to vertical mixing (called “salt stratification”). In addition, the thermocline, that is, the region of the ocean that features a rapid change in temperature with pressure, is visible in some temperature profiles.

## 2.2. GMM

GMM is a probabilistic approach for describing and classifying data. It attempts to fit (or “model” in the statistical sense) the data as a linear combination of multidimensional Gaussian distributions with unknown means and unknown standard deviations. Let  $\mathbf{X}$  be the array of  $N$  vertical profiles, each with  $D$  pressure levels, and let  $p(\mathbf{X})$  be the probability distribution function (PDF) representing the entire data set. GMM represents the PDF as a weighted sum of  $K$  Gaussian classes, indexed by  $k$ ; that is,

$$p(\mathbf{X}) = \sum_{k=1}^K \lambda_k \mathcal{N}(\mathbf{X}; \boldsymbol{\mu}_k, \Sigma_k). \quad (1)$$



**Figure 2.** Plot of 10% of the Argo temperature profiles, chosen at random, in the upper 1,000 dbar of the cleaned data set, along with the mean (solid line) and the mean plus or minus one standard deviation (dashed lines) across the entire data set.

Here  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  is the multidimensional Gaussian PDF with a vector of means  $\boldsymbol{\mu}_k$  and covariance matrix  $\boldsymbol{\Sigma}_k$ ; that is,

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right]}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}_k|}}. \quad (2)$$

The probability associated with class  $k$  is  $p(k) = \lambda_k$ . The probability of profile  $\mathbf{x}$  being in class  $k$  is  $p(k|\mathbf{x}) = \lambda_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) / p(\mathbf{x})$ , where the vector  $\mathbf{x}$  is a single profile taken from the complete array  $\mathbf{X}$  and  $p(\mathbf{x})$  is equation (1) with a single profile  $\mathbf{x}$  as the argument, that is, a normalizing factor. Both  $\mathbf{x}$  and  $\boldsymbol{\mu}_k$  are vectors of length  $D$ , and  $\boldsymbol{\Sigma}_k$  is a matrix of size  $D \times D$ .

Starting with random initial guesses for the classes, GMM proceeds by iteratively adjusting the means  $\boldsymbol{\mu}_k$  and standard deviations  $\boldsymbol{\Sigma}_k$  (i.e., the “parameters”) of the classes in order to maximize a logarithmic measure of likelihood; that is,

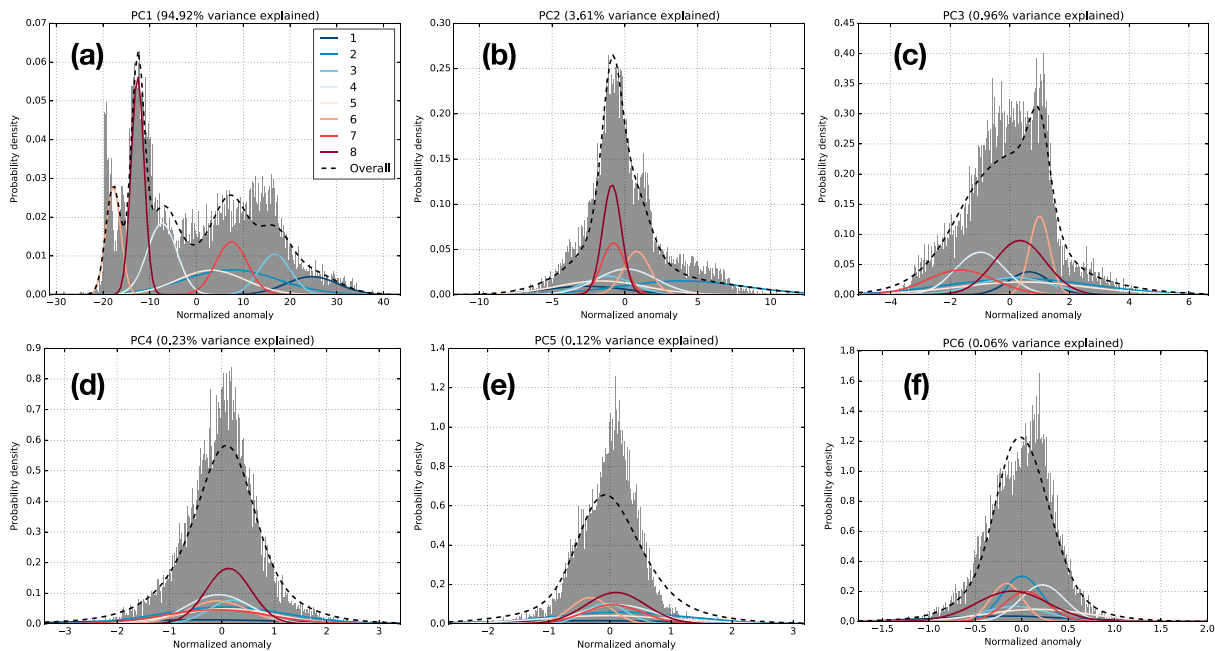
$$\log[p(\mathbf{X})] = \sum_{i=1}^N \log \left[ \sum_{k=1}^K \lambda_k \mathcal{N}(\mathbf{X}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]. \quad (3)$$

GMM uses an expectation-maximization approach, described in Maze et al. (2017). This algorithm monotonically converges on a local maximum. GMM is a generalization of  $k$ -means clustering, which only attempts to minimize in-group variance by shifting the means. By contrast, GMM attempts to identify means and standard deviations, allowing for some variation about the centers of the Gaussian distributions.

In our instance of GMM, each pressure level is treated as a “dimension,” and the Gaussian parameters are associated with each pressure level. However, we may not need all of these pressure levels to accurately describe the data set, as ocean temperature changes much more rapidly in the mixed layer and thermocline than in the interior. In order to reduce the computational complexity of the problem, we transform the profile data from pressure space to an alternative space using principal component analysis (PCA). Specifically, we calculate principal components (PCs) that capture a desired fraction of the vertical variability of the data set. Each eigenvector may be thought of as a “profile type” that describes a certain amount of variance in the data with pressure (note that this is not necessarily the same thing as a “typical profile”). We calculate  $J$  PCs via the transformation:

$$\mathbf{X}(z) = \sum_{j=1}^J \mathbf{P}(z, j) \mathbf{Y}(j), \quad (4)$$





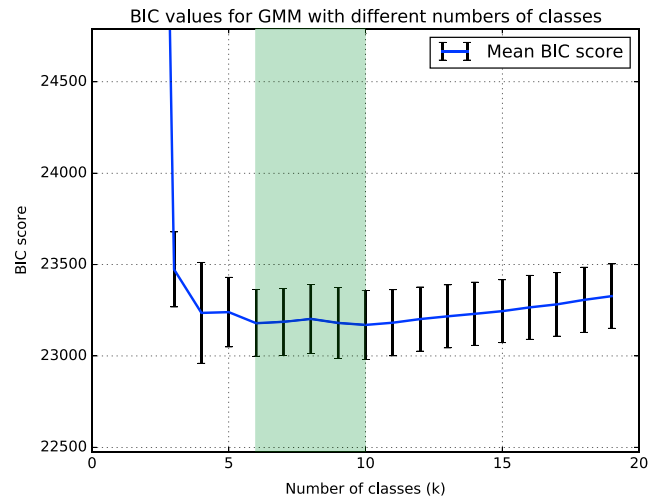
**Figure 3.** (a–f) Probability density functions for the (dimensionless) principal component amplitude coefficients associated with each profile, along with the Gaussian functions generated by Gaussian mixture modeling with  $K = 8$  classes. PC = principal component.

where  $z$  is the pressure level,  $J$  is the total number of PCs (index  $j$ ), and  $\mathbf{P}(z, j)$  is the transformation matrix between pressure space and PC space. This strategy is an example of “dimensionality reduction,” which is common in machine learning approaches.

We find that  $J = 6$  captures 99.9% of the variance in the vertical structure, which greatly reduces the number of dimensions needed to describe the Argo profile data used here, that is, from 194 pressure levels to 6 PCs. We refer to this data set as the “cleaned, compressed” data set. Nearly 95% of the variance is explained by the first PC (i.e., PC1), and the Gaussian functions associated with PC1 are relatively distinct, capturing the broad shape of the temperature distribution (Figure 3). For higher indexed PCs, the Gaussians overlap more, but their sum still makes up a representation of the temperature distribution that is sufficiently accurate for our purposes. The fact that we only need six PCs to capture 99.9% of the variance is consistent with the strong vertical coherence found in the SO, which is well described by an equivalent barotropic model (Karsten & Marshall, 2002). For more information on the PCs that we used in this work, see the supporting information (Figures S1 and S2).

We used a “training” data set, a subset of the cleaned, compressed data set, to estimate the parameters (i.e. the means and standard deviations) of the GMM classes. To generate the GMM training set, we randomly selected a single profile from each  $1^\circ \times 1^\circ$  bin. Each training data set contains 12,286 profiles (roughly 4% of the cleaned, compressed data set), distributed evenly in latitude/longitude space. Note that this subselection is not related cross-validation analysis, in which there are training and “test” data sets (Maze et al., 2017). Instead, we use a random subselection that is roughly uniform in latitude-longitude as our test data set, and then we apply the GMM model to the entire cleaned, compressed data set. As discussed in the supporting information, our results are not sensitive to our choice of test data set.

Once we have our test data set and calculate the optimized parameters (i.e., the means and standard deviations of the Gaussians), we then statistically represent (i.e., model) the entire cleaned, compressed data set with the fitted Gaussian model using optimized parameters. The end result is a probabilistic description of the cleaned, compressed Argo temperature profile data set in terms of a linear combination of Gaussian distributions that vary with pressure. Each profile then has a probability distribution across the classes, and the profile is assigned to the class with the highest probability. Our results are not sensitive to our choice of training data set (see supporting information Table S1).



**Figure 4.** Bayesian information criteria (BIC) scores versus the specified number of classes  $K$ . For each  $K$ , we calculate the BIC score 50 times using randomly selected profiles as discussed in the text. The means (solid blue line) and standard deviations (error bars) are shown for each  $K$ . The range of the smallest mean  $K$  values is indicated by green shading. GMM = Gaussian mixture modeling.

### 2.2.1. Selecting the Number of Classes

GMM does have one free parameter, that is, the maximum number of classes  $K$ . In order to determine the most appropriate value for  $K$ , we applied a statistical test, namely, a Bayesian information criterion (BIC). BIC uses an empirically formulated cost function that rewards likelihood and penalizes the number of classes  $K$ :

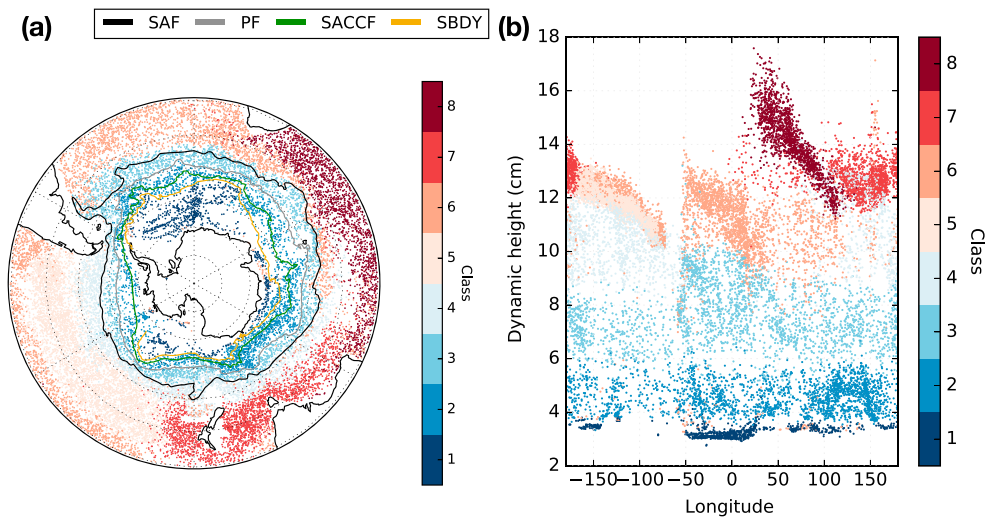
$$BIC(K) = -2\mathcal{L}(K) + N_f(K) \log(n), \quad (5)$$

where  $\mathcal{L}$  is a measure of likelihood,  $n$  is the number of profiles used in the BIC test, and  $N_f$  is the number of independent parameters to be estimated:

$$N_f(K) = K - 1 + KD + \frac{KD(D - 1)}{2}. \quad (6)$$

In this framework, the optimum value of  $K$  minimizes the BIC score. We perform a number of BIC tests, using different subsets of the data and different values of  $K$ , to estimate the distribution and variability of BIC. Using the roughly 300-km decorrelation scale of the SO as guidance (Ninove et al., 2016), we randomly select a profile from each  $4^\circ \times 4^\circ$  grid cell, returning 884 random profiles for each BIC test. We calculate BIC scores for each set of 884 random profiles (in PC space) using a range of classes  $K$  from 1 to 19 (Figure 4). For each value of  $K$ , we repeat the random selection and BIC process 50 times. BIC analysis does not feature a clear minimum, but instead, it suggests that the optimum value of  $K$  lies between 6 and 10.

It may seem counterintuitive that BIC does not return a single optimum value for  $K$ , but this is consistent with the nature of  $K$  as a weakly constrained free parameter that controls the level of complexity of the statistical description of the data set. Oceanography has a rich history of expertise-driven clustering using physical and biogeochemical criteria (e.g., potential vorticity minima and oxygen minima) and the fingerprints of various processes (e.g., gyre circulation). These descriptions can be arranged into hierarchies, from coarse/simple (e.g., two-layer quasi-geostrophic models) to rich and complex (e.g., the descriptions found in Talley, 2013). The level of detail required in the description depends on the application at hand. For example, a simple  $\beta$ -plane model is sufficient to explain the existence of gyres and western boundary currents; it constitutes a first-order description of gyres. Algorithmic clustering offers a robust way to traverse this hierarchy using a range of  $K$  values. Although statistical tests can be used as rough guides for choosing the number of classes, there is not necessarily a single ideal value for  $K$ . We explore the impact of  $K$  on our results in section 4.



**Figure 5.** (a) Gaussian mixture modeling-derived class distribution for  $K = 8$ , shown with four fronts of the Antarctic Circumpolar Current, that is, the Subantarctic Front (SAF), Polar Front (PF), Southern ACC Front (SACCF), and the southern boundary (SBDY; Kim & Orsi, 2014). (b) Class distribution shown in dynamic height space ( $\varphi_{1,500}^{300 \text{ dbar}}$ ). Note that only points with posterior probability  $\geq 0.9$  are shown. The classes are sorted by mean temperature, from coldest ( $k = 1$ ) to warmest ( $k = 8$ ).

### 3. Results

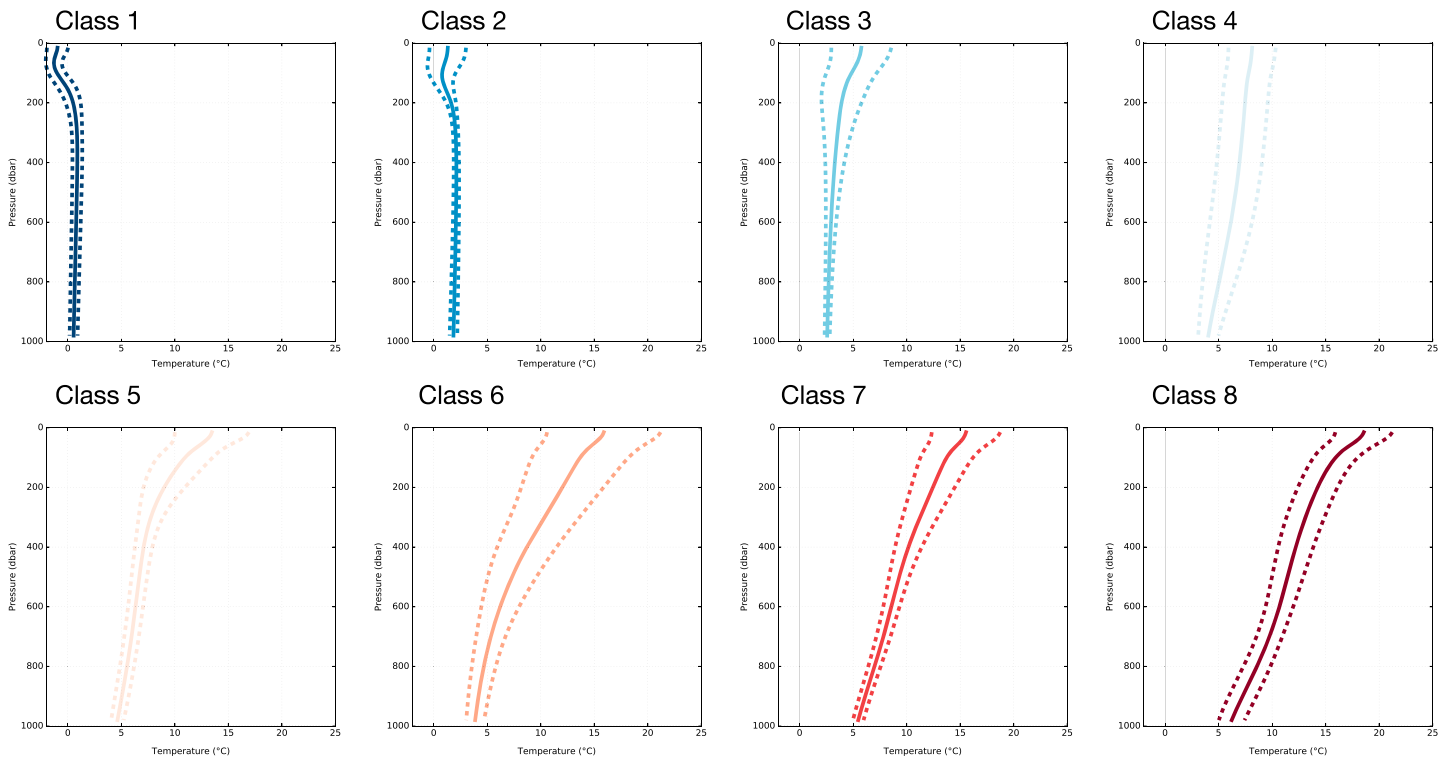
In order to identify patterns in the temperature structure of the SO, we describe the cleaned, compressed Argo temperature profile data set as a linear combination of multidimensional Gaussian functions that vary with pressure, using  $K = 8$  different classes. Despite the fact that GMM does not have access to the longitudes and latitudes of the profiles, it identifies spatially coherent structures, some of which are roughly demarcated by the fronts of the ACC as defined by Kim and Orsi (2014; Figure 5). For ease of interpretation, we sorted the classes by mean temperature (Table 1).

The class nearest Antarctica (class 1) extends throughout the Weddell Gyre and coastal Antarctica (Figure 5a). The mean temperature profile in this region is inverted; that is, it is colder near the surface and warmer in the interior (Figure 6). This near-Antarctic class coincides with regions of AABW export (Ohshima et al., 2013; Orsi et al., 1999), the subpolar Weddell and Ross Gyres, and its northern extent approximately corresponds with the classical SBDY of the ACC (Kim & Orsi, 2014). This class occupies a narrow range in dynamic height space, with a class mean and standard deviation of  $3.3 \pm 0.2 \text{ cm}$  ( $\varphi_{1,500}^{300 \text{ dbar}}$ ; Figure 5b); it is fairly distinct from the other classes; that is, class 1 profiles are rarely found north of the SBDY. For reference, Kim and Orsi (2014) associate the SBDY with the 3.1-cm dynamic height contour ( $\varphi_{1,500}^{500 \text{ dbar}}$ ). As their limits of integration over pressure are different than ours, this value of dynamic height

**Table 1**  
Temperature Statistics for Each Class, Using Values From Every Pressure Level

Class	Number of profiles	Mean	Standard deviation	Minimum	Maximum
1	10,680	0.48	0.81	-2.11	2.52
2	33,031	1.83	0.72	-1.87	8.89
3	40,268	3.38	1.50	-1.82	19.70
4	39,619	6.36	2.24	-1.85	17.17
5	48,252	7.32	2.56	2.76	25.37
6	48,770	8.22	4.49	-1.88	27.56
7	38,682	9.70	3.07	3.25	27.11
8	25,130	11.57	3.43	3.56	28.08

*Note.* All temperature statistics are shown in degrees Celsius. The classes have been sorted by mean temperature, calculated using values from all pressure levels.



**Figure 6.** Temperature profile statistics, separated by class, as functions of pressure. Shown are the mean (solid lines) and the mean plus or minus one standard deviation (dashed lines) for all profiles in the indicated class.

is not directly applicable to our data, but it is roughly consistent with the gap between classes 1 and 2 in our analysis (Figure 5b). Assuming that the data feature sufficiently uniform spatial coverage, gaps in dynamic height space may be indicative of fronts, as they may suggest sharp gradients in dynamic height over relatively short physical distances. We do not pursue this analysis further here. For an in-depth analysis of SO front positions, see Sokolov and Rintoul (2009), for example.

The second coldest class (class 2) is a circumpolar class with profiles that sit north of the SBDY and south of the PF across all longitudes; it is the dominant class in the dynamic height range 4–6 cm, with a class mean value of  $4.8 \pm 0.7$  cm ( $\phi_{1,500}^{300 \text{ dbar}}$ , Figure 5). Its mean profile is also inverted, though not as sharply as the mean profile of class 1 (Figure 6). A second circumpolar class (class 3) sits roughly north of the PF and south of the SAF. In dynamic height space, class 3 is found between roughly 6–8 cm, except in the Atlantic sector, where it extends to roughly 10 cm. For reference, Kim and Orsi (2014) associate the PF with the 5.0-cm dynamic height contour and the SAF with the 7.0-cm dynamic height contour ( $\phi_{1,500}^{500 \text{ dbar}}$ ). These values are roughly consistent with (but not directly comparable to) the gap positions in our data. Unlike the first two classes, the mean profile of class 3 is not inverted; that is, it gets colder with pressure. The presence of these two circumpolar classes is consistent with the homogenizing influence of the ACC, which typically encourages mixing along the strong jets associated with fronts and suppresses mixing across them (Ferrari & Nikurashin, 2010).

The profiles assigned to class 4 are mostly located north of the SAF in the Pacific and Indian sectors, roughly coinciding with regions of Subantarctic Mode Water (SAMW) and Antarctic Intermediate Water (AAIW) formation in the Pacific Ocean and south of Australia (Sallée et al., 2010b). Despite its relatively narrow range in latitude, class 4 profiles occupy a broad, distinct range in dynamic height space in the Pacific sector, with a class mean of  $11 \pm 1.5$  cm. The mean vertical profile associated with class 4 changes relatively gently with pressure, with no clear thermocline and a relatively large standard deviation across all pressures.

Profiles assigned to class 5 are mostly found in the Pacific sector, in a region associated with the export of SAMW and AAIW from the surface ocean into the interior thermocline (Iudicone et al., 2007; Jones et al., 2016). In contrast with class 4, class 5 occupies a relatively large range in latitude and a relatively small range

**Table 2**  
*Posterior Probabilities for Each Class, Divided Into Four Unequal Intervals*

Class	[0.0, 0.50)	[0.50, 0.75)	[0.75, 0.9)	[0.9, 1.0]
1	<1	4	4	91
2	<1	11	11	77
3	<1	14	16	70
4	1	18	20	61
5	<1	7	8	84
6	<1	9	8	82
7	<1	19	17	64
8	<1	13	12	75

*Note.* Each row shows the percentage of profiles assigned to that class with posterior probabilities in the indicated range.

in dynamic height, with a mean and standard deviation of  $12 \pm 0.7$  cm. The mean vertical profile has a clear thermocline over the upper 400 dbar of the ocean, with a standard deviation that narrows considerably with pressure. This class spatially coincides with the southern part of the South Pacific gyre, suggesting that gyre circulation tends to homogenize properties in this region.

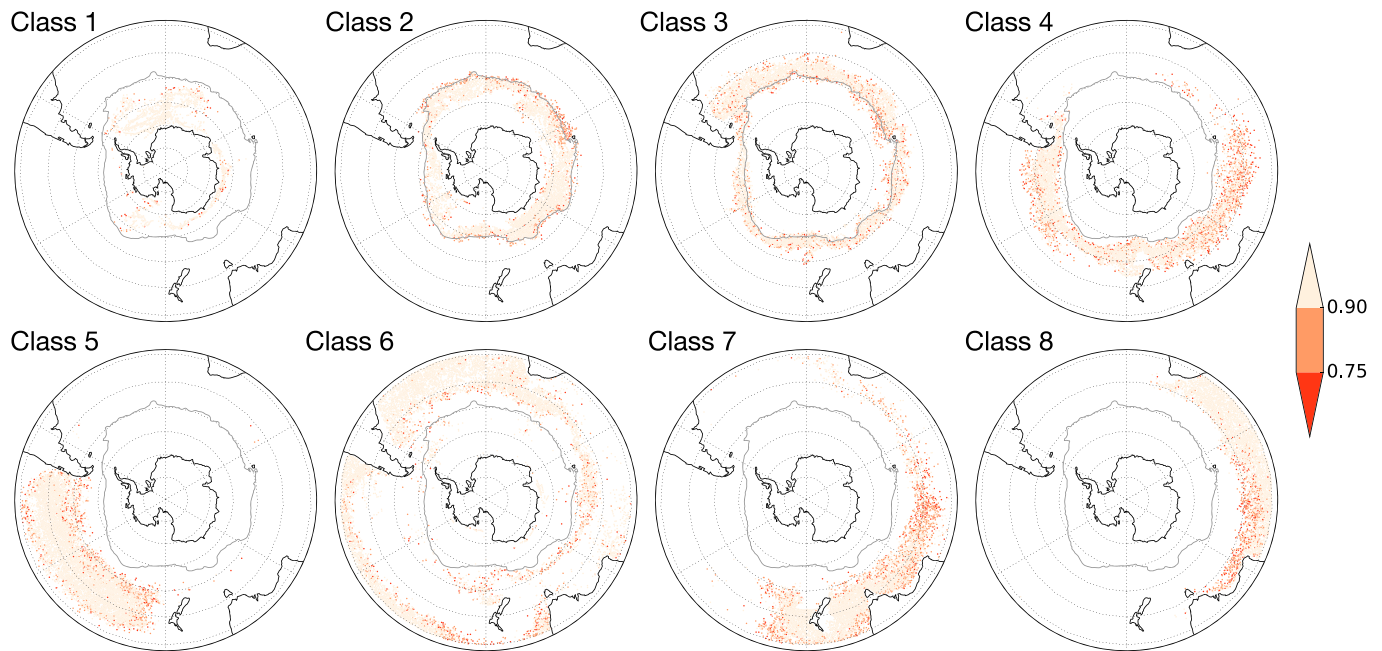
Class 6 highlights warmer subtropical waters and is mostly found in the Atlantic and Pacific sectors; it partially extends into the Indian sector, where it sits just north of the SAF. From the surface to well into the interior, class 6 features some of the largest standard deviations of any class, suggesting that class 6 consists of a wide variety of profiles; it can potentially be split into a number of smaller classes. Classes 7 and 8 are also warmer subtropical classes, with class 7 found mostly near Australia and New Zealand and class 8 found almost exclusively in the Indian sector. Much of class 8 spatially coincides with the Indian Ocean gyre. The spatial extent of class 8 near South Africa suggests that the Agulhas Current influences the temperature structure in that region. The mean vertical profiles of classes 7 and 8 are similar, although class 7 features higher variability near the surface and class 8 features slightly warmer surface temperatures. The higher variability in class 7 may be due to the overlap of profiles in this class with a wider range of surface current features (e.g., boundary currents around Australia and New Zealand, whereas class 8 largely overlaps with the Indian Ocean gyre.

For a selected temperature profile, GMM predicts the probability distribution across all  $K$  classes. That is, it calculates the probabilities that the profile belongs to each class  $k$ . Next, the algorithm assigns the profile to the class with the highest probability. Note that the sum of the posterior probabilities across all classes is one. Since these probabilities are calculated with the full data set available, they are referred to as posterior probabilities. The posterior probabilities are useful in their own right, as measures of confidence in GMM's assignment of a profile to a particular class.

For our implementation of GMM on Argo temperature data, over 86% of the class assignments have posterior probabilities greater than 0.75, and over 74% of all class assignments have posterior probabilities greater than 0.9 (Table 2). Class 1 features an especially high percentage of very high posterior probabilities; over 90% of assignments into class 1 have posterior probabilities greater than or equal to 0.9. Outside of the Weddell Gyre, we find the lowest posterior values in the Ross Sea and a few near-coastal areas (Figure 7). The low posterior values could possibly be due to seasonal variability that is not well represented by a single class. Classes 2 and 3 also feature high posterior probabilities, for which over 70% of assignments have values greater than or equal to 0.9. For both of these classes, we find relatively low posterior probabilities upstream of Kerguelen Island (KI), clustered around the PF. The area around KI is affected by upwelling, mixing, and the confluence of the Agulhas Retroflexion and the ACC (Sallée et al., 2010b), and it also features relatively high eddy diffusivities (Klocker & Abernathy, 2014). The profiles in that area are influenced by a number of competing processes and may be difficult to unambiguously separate into clear groups when using a value of  $K$  appropriate for the entire SO.

Although over 60% class 4 profiles have posterior values greater than or equal to 0.9, class 4 features some relatively low posterior values compared with the other classes, especially in the Indian sector north of the SAF. In the Pacific sector, we find relatively low posteriors along the boundary between classes 4 and 5.





**Figure 7.** Posterior probabilities for each class assignment, given the full cleaned, compressed data set, shown together with the PF for reference (Kim & Orsi, 2014).

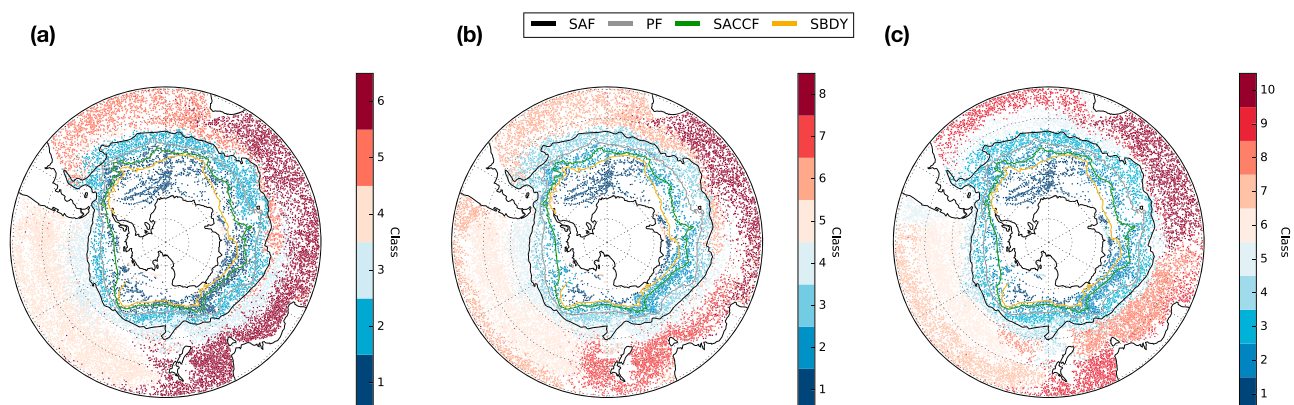
Class 5 has a core of profiles with posterior values greater than or equal to 0.9, with lower values all along its boundary. We find similar patterns for classes 6–8, except in the Indian sector between 60–120°E, north of the SAF. This region, which is downstream of Kerguelen Plateau, is characterized by relatively low posterior values for classes 4, 7, and 8. In general, although GMM performs well in all ocean basins, in terms of clear class separation with high posterior probabilities, its performance is somewhat weaker in the Indian sector.

#### 4. Discussion

Here we explore the sensitivity of our results to the maximum number of classes  $K$ . We also explore a possible alternative to PCA that may be useful for incorporating salinity into our analysis, namely, functional PCA.

##### 4.1. Sensitivity to Number of Classes $K$

In section 2, we estimated that the optimum number of classes  $K$  lies between 6 and 10. The weak constraint suggested by BIC allows for some tuning depending on the desired level of complexity in the description of the data set. Using  $K = 6$  classes is sufficient to capture most of the large-scale structures identified in the  $K = 8$  case, but there are some significant differences (Figures 8a and 8b). Specifically, there is one fewer



**Figure 8.** Comparison of Gaussian mixture modeling-derived classes, shown for (a) 6 classes, (b) 8 classes, and (c) 10 classes, along with fronts of the Antarctic Circumpolar Current (Kim & Orsi, 2014). SAF = Subantarctic Front; PF = Polar Front; SACCF = Southern ACC Front; SBDY = Southern Boundary.

circumpolar class, as classes 1–3 are reduced to classes 1–2 that sit roughly on either side of the PF. In the Pacific sector, classes 4 and 5 merge into the new class 4. In the Indian sector, classes 7–8 merge into the new class 6 that sits north of the SAF and south of Australia. We see that the overall description of ocean structure is simpler with  $K = 6$ ; it is still a physically reasonable description of ocean temperature structure, with circumpolar classes and clusters that span the major basins, but it lacks some of the subtleties found in the  $K = 8$  map.

As expected, the  $K = 10$  case features more structure than the  $K = 8$  case, and it is still a physically reasonable distribution (Figures 8b and 8c). Classes 1–3 are still near Antarctic or circumpolar classes; the additional structure all appears north of the SAF. In the Pacific basin, the boundary between the  $K = 8$  classes 5 and 6 and the  $K = 10$  classes 6 and 7 is shifted poleward, and a new class 5 is found along the eastern Pacific, along the South American coast. The  $K = 10$  class 8 is found south of Australia, which in the  $K = 8$  class is not a distinct class. Interestingly, in the  $K = 10$  case we find more profiles above 0.9 posterior probability in the Indian sector, specifically in the region north of the SAF and between the longitudes of 60–120°E. Increasing  $K$  allowed for a more likely set of class assignments in this previously troublesome region. So regions of low posterior probabilities may suggest a need for a higher value of  $K$ .

#### 4.2. Including Additional Variables

In this work, we define classes using temperature profiles. A more general description of ocean structure may include some combination of additional fields (e.g., salinity, density, potential vorticity, and biogeochemical variables). Including additional variables in the GMM analysis is not necessarily trivial, as there are different approaches, and each approach has advantages and limitations that need to be thoroughly evaluated. Perhaps, the simplest approach is to standardize the additional variables in the same way as the temperature fields, such that each field is expressed in terms of standard deviations relative to the mean at each depth level. This approach has the potential disadvantage that a  $1\sigma$  variation in temperature has the same impact on class structure as a  $1\sigma$  variation in salinity, which is not necessarily physically realistic. Statistically similar variations in temperature and salinity need not have similar impacts on ocean structure. Another approach is to scale by parameters from the linear equation of state (EOS), namely, the thermal expansion coefficient  $\alpha$  for temperature and the coefficient of haline contraction  $\beta$  for salinity. Using this method, variations in temperature and salinity would impact density in a manner that is physically constrained by the EOS, that is,  $\rho = \rho_0 [1 - \alpha(T - T_0) + \beta(S - S_0)]$ . However, this approach would only be valid in the neighborhood of the reference values  $(T_0, S_0)$  in  $T$ - $S$  space, over which the linear EOS is a good approximation of the full, nonlinear EOS. This limitation would likely be problematic in the SO, where nonlinear terms in the EOS play an important role in the formation and layering of AAIW and AABW (Nycander et al., 2015). Still, another approach would be to classify profiles based on density, but there are a number of different approaches to defining density that would need to be treated with care (e.g.,  $\sigma_0$ ,  $\sigma_1$ , and neutral density  $\gamma_n$ ).

We used PCA to reduce the dimensionality of our Argo temperature profile data set. An alternative approach is to use functional principal component analysis (fPCA), in which PCA is performed on functions instead of the original data. In Pauthenet et al. (2017), the authors represent vertical temperature and salinity profiles from the Southern Ocean State Estimate (Mazloff et al., 2010) as linear combinations of B-spline basis functions and apply fPCA to the resulting spline functions. They use the PCs to examine large-scale structures such as fronts in the SO. Their approach offers another objective way to define structural boundaries and could be used in concert with the GMM approach outlined in this work. This could offer yet another possible way to introduce salinity into the GMM analysis, which is especially relevant for stratification south of the PF (Pollard et al., 2002).

## 5. Conclusions

We applied GMM, an unsupervised classification scheme, to SO Argo temperature data above 1,000 dbar. Without using longitude or latitude information, GMM identified spatially coherent patterns in the vertical temperature structure. The GMM-derived classes broadly coincide with large-scale circulation and stratification features, including regions of AABW formation and upwelling (i.e., adjacent to Antarctica), the ACC, formation and export pathways of SAMW and AAIW, subtropical gyre circulation, and the Agulhas Current and associated retroflexion. We may say that GMM identifies *domains* in oceanographic data, including gyre-dominated domains and circumpolar domains. GMM can be used to define these domains in a method that respects the structure of the data, as opposed the simpler but physically unrealistic process of defin-

ing domains by simply drawing rectangular boxes in latitude-longitude space. GMM also makes use of the interior structure of the data, as opposed to only using surface variables like sea surface height. The class boundaries broadly coincide with several classically defined fronts of the ACC, and the circumpolar classes mostly occupy distinct regions in dynamic height space, indicating that GMM has identified physically distinct profile types using only vertical temperature data. High posterior probability distributions indicate regions where the classes are distinct and statistically separate, whereas regions with low posterior probability indicate boundaries between classes and/or regions of mixing influenced by a number of different temperature structures. GMM may offer an alternative, complementary method for describing SO spatial variability, and it is potentially useful for objectively and automatically comparing structures across different observational and modeling data sets.

### Acronyms

AABW—Antarctic Bottom Water  
 AAIW—Antarctic Intermediate Water  
 ACC—Antarctic Circumpolar Current  
 BIC—Bayesian information criterion  
 fPCA—Functional principal component analysis  
 GDAC—Global Data Assembly Center  
 GMM—Gaussian mixture modeling  
 PC—Principal component  
 PCA—Principal component analysis  
 PDF—Probability distribution function  
 SAMW—Subantarctic Mode Water

### Acknowledgments

This study is supported by grants from the Natural Environment Research Council (NERC), including (1) The North Atlantic Climate System Integrated Study (ACSIS) (grant NE/N018028/1 [authors D. J. and E. S.]) and (2) Ocean Regulation of Climate by Heat and Carbon Sequestration and Transports (ORCHESTRA) (grant NE/N018095/1 [authors E. S. and A. M.]). H. H. was funded by a NERC DTP Research Experience Placement over the summer of 2017 (grant NE/L002434/1). Argo floats data were collected and made freely available by the International Argo Program and the national programs that contribute to it (<http://www.argo.ucsd.edu> and <http://argo.jcommops.org>). The Argo Program is part of the Global Ocean Observing System. Argo floats data and metadata are available from the Global Data Assembly Centre (Argo GDAC), <http://doi.org/10.17882/42182>. The analysis software used in this manuscript was written using Python and the *scikit-learn* machine learning library (<http://scikit-learn.org/stable/>). Version 1.0 of the scripts we used are available via github (<https://github.com/DanJonesOcean/OceanClustering/releases/tag/v1.0>). D. J. thanks Chris Lowder for python support. We are grateful to Y. S. Kim for providing us with Southern Ocean front position data. Finally, we thank Guillaume Maze and three anonymous reviewers, whose feedback greatly improved the quality of our work.

### References

- Emery, W. J. (2003). Water types and water masses. In J. R. Holton, J. A. Curry, & J. A. Pyle (Eds.), *Encyclopedia of atmospheric sciences* (pp. 1556–1567). New York: Elsevier. <https://doi.org/10.1016/b0-12-227090-8/00279-7>
- Ferrari, R., & Nikurashin, M. (2010). Suppression of eddy diffusivity across jets in the Southern Ocean. *Journal of Physical Oceanography*, *40*, 1501–1519. <https://doi.org/10.1175/2010JPO4278.1>
- Fletcher, S. M., Gruber, N., Jacobson, A. R., Doney, S. C., Dutkiewicz, S., Gerber, M., et al. (2006). Inverse estimates of anthropogenic CO<sub>2</sub> uptake, transport, and storage by the ocean. *Global Biogeochemical Cycles*, *20*, GB2002. <https://doi.org/10.1029/2005gb002530>
- Frölicher, T. L., Sarmiento, J. L., Paynter, D. J., Dunne, J. P., Krasting, J. P., & Winton, M. (2015). Dominance of the Southern Ocean in anthropogenic carbon and heat uptake in CMIP5 models. *Journal of Climate*, *28*(2), 862–886. <https://doi.org/10.1175/jcli-d-14-00117.1>
- Hanawa, K., & Talley, L. (2001). Mode waters. In G. Siedler & J. Church (Eds.), *Ocean circulation and climate, International Geophysics Series* (pp. 373–386). Cambridge, MA, USA: Academic Press.
- Herraiz-Borreguero, L., & Rintoul, S. R. (2011). Subantarctic mode water: Distribution and circulation. *Ocean Dynamics*, *61*(1), 103–126. <https://doi.org/10.1007/s10236-010-0352-9>
- Holt, H. J., & Jones, D. C. (2018). Southern Ocean temperature profile classification tool. <https://doi.org/10.5281/zenodo.1543106>
- Iudicone, D., Rodgers, K., Schopp, R., & Madec, G. (2007). An exchange window for the injection of Antarctic Intermediate Water into the South Pacific. *Journal of Physical Oceanography*, *37*, 31–49. <https://doi.org/10.1175/JPO2985.1>
- Jones, D. C., Meijers, A. J. S., Shuckburgh, E., Sallée, J.-B., Haynes, P., McAufield, E. K., & Mazloff, M. R. (2016). How does Subantarctic Mode Water ventilate the Southern Hemisphere subtropics? *Journal of Geophysical Research: Oceans*, *121*, 6558–6582. <https://doi.org/10.1002/2016jc011680>
- Karsten, R. H., & Marshall, J. (2002). Constructing the residual circulation of the ACC from observations. *Journal of Physical Oceanography*, *32*, 3315–3327. [https://doi.org/10.1175/1520-0485\(2002\)032<3315:CTRCOT>2.0.CO;2](https://doi.org/10.1175/1520-0485(2002)032<3315:CTRCOT>2.0.CO;2)
- Kim, Y. S., & Orsi, A. H. (2014). On the variability of Antarctic Circumpolar Current fronts inferred from 1992–2011 altimetry\*. *Journal of Physical Oceanography*, *44*(12), 3054–3071. <https://doi.org/10.1175/JPO-D-13-0217.1>
- Klocker, A., & Abernathy, R. (2014). Global patterns of mesoscale eddy properties and diffusivities. *Journal of Physical Oceanography*, *44*(3), 1030–1046. <https://doi.org/10.1175/jpo-d-13-0159.1>
- Lumpkin, R., & Speer, K. (2007). Global ocean meridional overturning. *Journal of Physical Oceanography*, *37*, 2550–2562. <https://doi.org/10.1175/JPO3130.1>
- Maze, G., Mercier, H., Fablet, R., Tandeo, P., Radcenco, M. L., Lenca, P., et al. (2017). Coherent heat patterns revealed by unsupervised classification of Argo temperature profiles in the North Atlantic Ocean. *Progress in Oceanography*, *151*, 275–292. <https://doi.org/10.1016/j.pocean.2016.12.008>
- Mazloff, M. R., Heimbach, P., & Wunsch, C. (2010). An eddy-permitting Southern Ocean state estimate. *Journal of Physical Oceanography*, *40*(5), 880–899. <https://doi.org/10.1175/2009jpo4236.1>
- Naveira-Garabato, A. C. N., Ferrari, R., & Polzin, K. L. (2011). Eddy stirring in the Southern Ocean. *Journal of Geophysical Research*, *116*, C09019. <https://doi.org/10.1029/2010jc006818>
- Naveira-Garabato, A. C. N., Jullion, L., Stevens, D. P., Heywood, K. J., & King, B. A. (2009). Variability of Subantarctic Mode Water and Antarctic Intermediate Water in the Drake Passage during the late-twentieth and early-twenty-first centuries. *Journal of Climate*, *22*(13), 3661–3688. <https://doi.org/10.1175/2009jcli2621.1>

- Ninove, F., Le Traon, P. Y., Remy, E., & Guinehut, S. (2016). Spatial scales of temperature and salinity variability estimated from Argo observations. *Ocean Science*, *12*(1), 1–7. <https://doi.org/10.5194/os-12-1-2016>
- Nycander, J., Hieronymus, M., & Roquet, F. (2015). The nonlinear equation of state of sea water and the global water mass distribution. *Geophysical Research Letters*, *42*, 7714–7721. <https://doi.org/10.1002/2015GL065525>
- Ohshima, K. I., Fukamachi, Y., Williams, G. D., Nihashi, S., Roquet, F., Kitade, Y., et al. (2013). Antarctic Bottom Water production by intense sea-ice formation in the Cape Darnley polynya. *Nature Geoscience*, *6*(3), 235–240. <https://doi.org/10.1038/ngeo1738>
- Orsi, A. H., Johnson, G. C., & Bullister, J. L. (1999). Circulation, mixing, and production of Antarctic Bottom Water. *Progress in Oceanography*, *43*(1), 55–109. [https://doi.org/10.1016/S0079-6611\(99\)00004-X](https://doi.org/10.1016/S0079-6611(99)00004-X)
- Orsi, A., Whitworth, T., & Nowlin, W. (1995). On the meridional extent and fronts of the Antarctic Circumpolar Current. *Deep Sea Research Part I*, *42*(5), 641–673.
- Pauthenet, É., Roquet, F., Madec, G., & Nerini, D. (2017). A linear decomposition of the Southern Ocean thermohaline structure. *Journal of Physical Oceanography*, *47*, 29–47. <https://doi.org/10.1175/JPO-D-16-0083.s1>
- Pollard, R. T., Lucas, M. I., & Read, J. F. (2002). Physical controls on biogeochemical zonation in the Southern Ocean. *Deep Sea Research Part II*, *49*(16), 3289–3305. [https://doi.org/10.1016/S0967-0645\(02\)00084-X](https://doi.org/10.1016/S0967-0645(02)00084-X)
- Purkey, S. G., & Johnson, G. C. (2010). Warming of global abyssal and deep Southern Ocean Waters between the 1990s and 2000s: Contributions to global heat and sea level rise budgets\*. *Journal of Climate*, *23*(23), 6336–6351. <https://doi.org/10.1175/2010jcli3682.1>
- Sallée, J., Morrow, R., & Speer, K. (2008). Eddy heat diffusion and Subantarctic Mode Water formation. *Geophysical Research Letters*, *35*, L05607. <https://doi.org/10.1029/2007GL032827>
- Sallée, J., Shuckburgh, E., Bruneau, N., Meijers, A., Bracegirdle, T., Wang, Z., & Roy, T. (2013). Assessment of Southern Ocean water mass circulation and characteristics in CMIP5 models: Historical bias and forcing response. *Journal of Research: Oceans*, *118*, 1830–1844. <https://doi.org/10.1002/jgrc.20135>
- Sallée, J., Speer, K., Rintoul, S., & Wijffels, S. (2010a). Southern Ocean thermocline ventilation. *Journal of Physical Oceanography*, *40*, 509–529. <https://doi.org/10.1175/2009JPO4291.1>
- Sallée, J.-B., Speer, K., Rintoul, S., & Wijffels, S. (2010b). Southern Ocean thermocline ventilation. *Journal of Physical Oceanography*, *40*(3), 509–529. <https://doi.org/10.1175/2009jpo4291.1>
- Sokolov, S., & Rintoul, S. R. (2009). Circumpolar structure and distribution of the Antarctic Circumpolar Current fronts: 1. Mean circumpolar paths. *Journal of Geophysical Research*, *114*, C11018. <https://doi.org/10.1029/2008JC005108>
- Talley, L. (2013). Closure of the global overturning circulation through the Indian, Pacific, and Southern Oceans: Schematics and transports. *Oceanography*, *26*(1), 80–97. <https://doi.org/10.5670/oceanog.2013.07>