

WHAT ARE THE ECOREGIONS OF THE OCEAN?

By Matias I. Bofarull Oddo

The University of British Columbia
Earth, Ocean, and Atmospheric Sciences

ABSTRACT	2
1. UNSUPERVISED CLUSTERS	3
1.1 OBTAINING DATASETS	3
1.2 PRE-PROCESSING	3
1.3 CAPTURING VARIANCE	7
1.3.1 PC 1	8
1.3.2 PC 2	9
1.3.3 PC 3	10
1.4 EXPLORING VARIANCE	10
1.5 HIERARCHICAL CLUSTERING	14
2. OCEANIC ECOREGIONS	15
2.1 WORLD BANDS	15
2.2 CLUSTERS WITHIN CLUSTERS	19
2.3 BOUNDARY STRENGTH	21
5. REFERENCES	23
APPENDIX	24
FINAL PAGE	53

Important: This report features graphics with black background.
Please abstain from printing this report.

ABSTRACT

As terrestrial animals we know of land ecoregions like grasslands, deserts, tundras, and rainforests. It is the variability of nutrients and energy that differentiates one ecoregion from another. But what are the ecoregions of the ocean? Regrettably, the most widely used oceanic partitions, like EEZs, are not based in ecologically-meaningful criteria. In order to study and steward marine resources, the boundaries between ecoregions should be congruent with an ecosystem's ability to support life, which is in turn determined by the energetic and nutritional signature of a region patch of seawater. This present study successfully uses unsupervised machine learning on World Ocean sea surface data, specifically Principal Component Analysis (PCA) and hierarchical clustering, to identify the location and strength of emergent boundaries between oceanic ecoregions.

1. UNSUPERVISED CLUSTERS

1.1 OBTAINING DATASETS

This study uses monthly climatology NOAA datasets from 2013 World Ocean Atlas (WOA2013). There are four ecologically-relevant factors related to the euphotic inorganic stock required for marine photosynthetic primary productivity: Sea Surface Nitrate, Sea Surface Silicate, Sea Surface Phosphate, and Sea Surface Oxygen. In addition, there are two factors related to seawater properties: Sea Surface Temperature and Sea Surface Salinity. WOA2013 offers these datasets with a spatial resolution of 1° of latitude/longitude, yielding world map matrices 180 rows tall (latitude 90N to 90S) by 360 columns wide (longitude 179E to 179W). Nitrate, Phosphate, Silicate, and Oxygen are measured in micromoles per Liter, Temperature is measured in degrees Celsius, and Salinity is measured in Practical Salinity Unit (PSU).

These six factors describe the chemical properties and inorganic stock of seawater, and the interplay of these variables has ecological significance. For example, temperature and salinity are used to calculate seawater density, which in turn is used to calculate the mixed layer depth, or the depth at which nutrients like nitrate, silicate, and phosphate are suspended for phytoplankton to use, which are in turn grazed by zooplankton that need oxygen for respiration.

1.2 PRE-PROCESSING

Extensive pre-processing is required to prepare the WOA2013 datasets for PCA and clustering. This pre-processing includes removing gaps in data from NaNs, detecting what are the latitude and longitude coordinates of oceanic versus landmass pixels, and numerically convert absolute magnitude values into relative variance for each variables. Before starting pre-processing, an exploration of the data using maps and histograms help visualize the seasonal variability within each factor. To get an idea of the distribution of raw values, the overall maxima, minima, and mean are calculated for each factor. This exploration is attached in the Appendix, specifically Supplemental Figures 1 to 12.

For each of the six factors (nitrate, silicate, phosphate, oxygen, temperature, and salinity) there is a monthly snapshot of the spatial variability during a given month, for a total of twelve seasonal matrices. These twelve snapshots are then stacked to create a 3D matrix per oceanic factor, for a total of six 3D matrices. Figure 1 illustrates the process to create these multifactorial matrices.

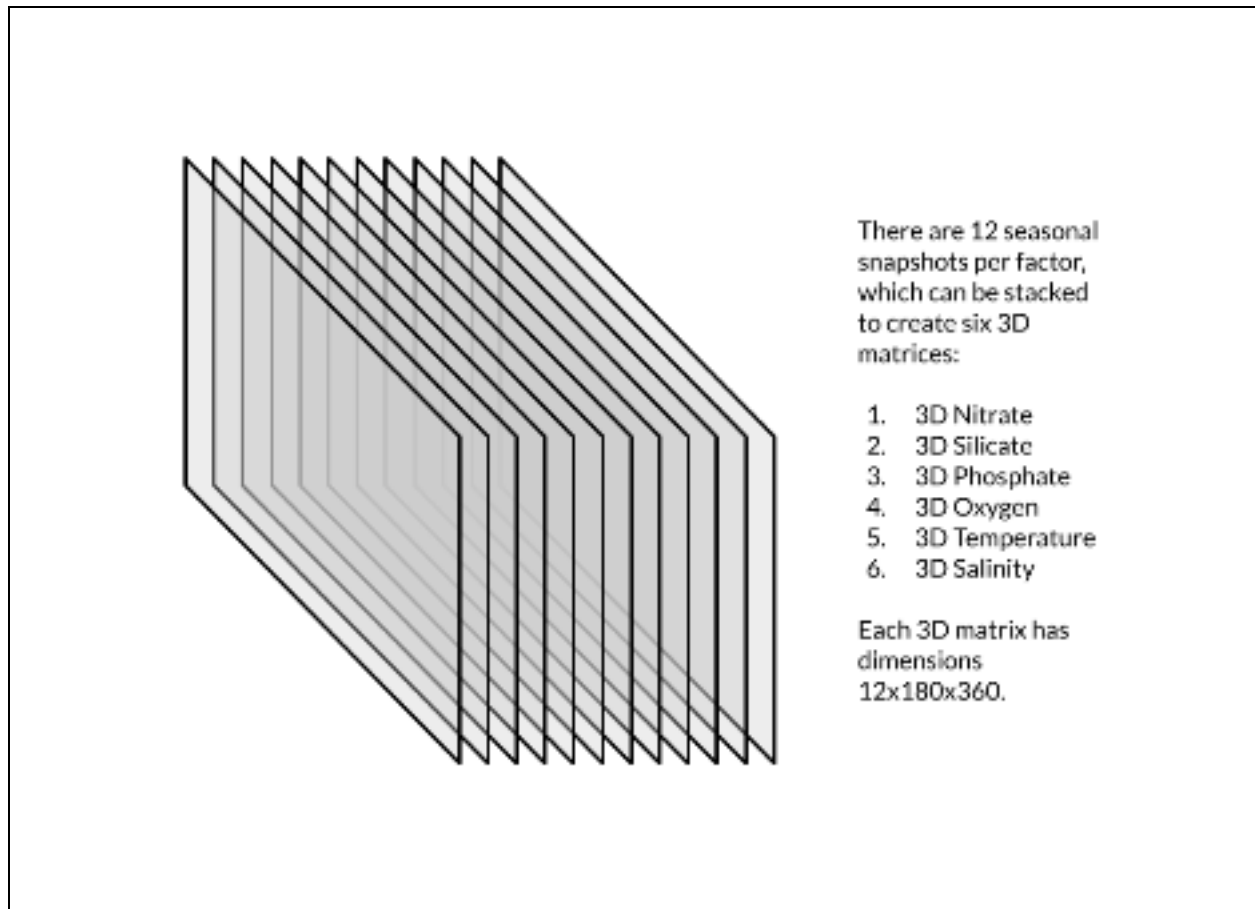


Figure 1. Seasonal snapshots stacked to create 3D matrices, one per oceanic factor.

Because seasonal measurements are inconsistent, some oceanic pixels have no data. In order to preserve seasonal variability without compromising analysis with unwanted NaNs, the maximum, minimum, and mean values per oceanic pixel were extracted per factor (Figure 2). Luckily, across all 3D matrices, in no case a single oceanic pixel happened to be NaN for twelve consecutive months, only landmass pixels exhibited this problematic property. This fortunate development resulted in all extracted matrices to have full pixel coverage for the maximum, minimum, and mean values for that variable.

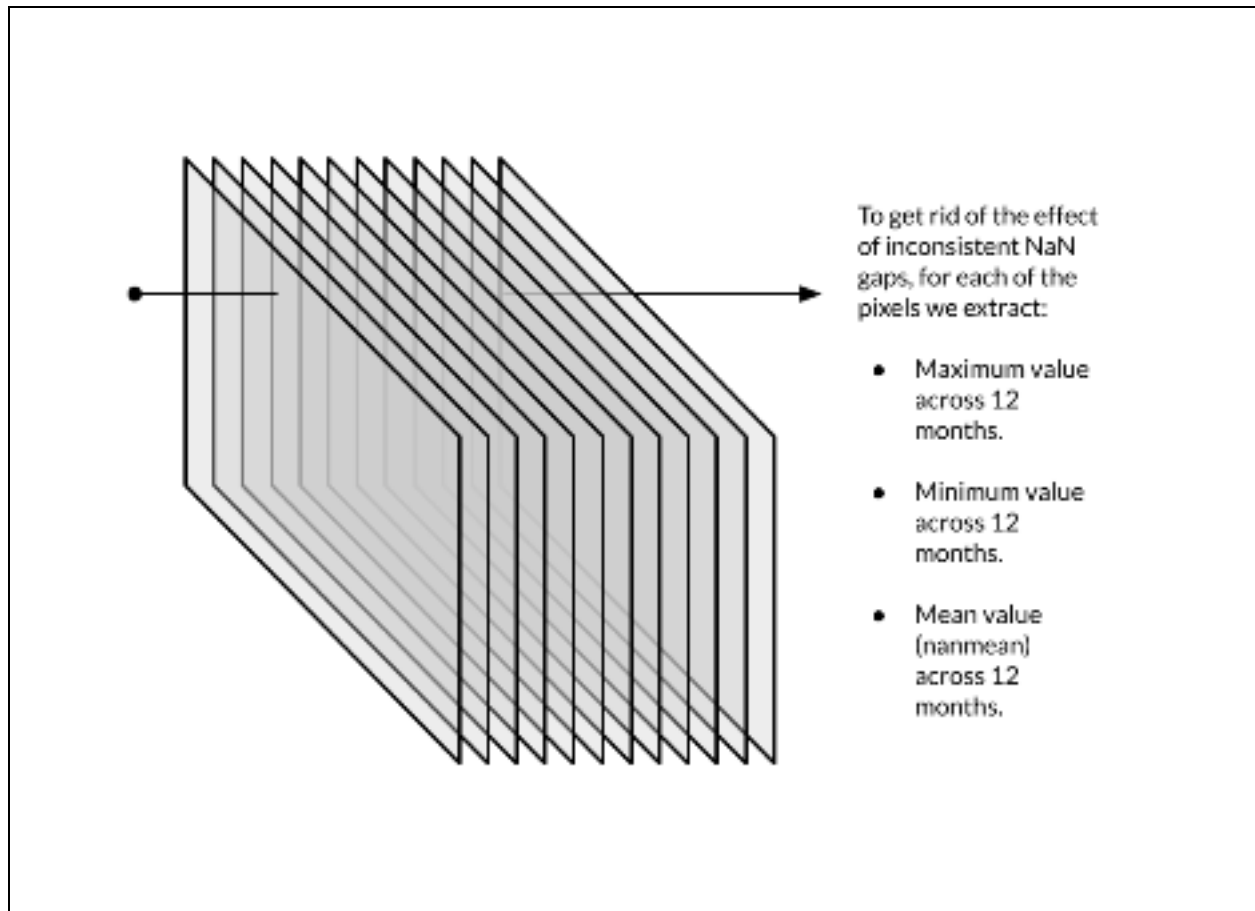


Figure 2. Extraction of values per oceanic pixel per factor to create the 18 variable matrices.

Extracting the maximum, minimum, and mean values for each pixel for each of the six factors results in six matrices for maximum values, six matrices for minimum values, and six matrices for mean values, for a total of eighteen variables in 2D matrix form. Because PCA is not concerned with the magnitude of values, but with the relative variability across values, all oceanic pixels within a variable were rescaled using the `rescale` function. This results in the internal variance within each of the eighteen variables to be numerically captured between 0 and 1, where the lowest value of the array becomes 0 and the largest 1, and the rest populates the intermediate range as decimals. The resulting eighteen variables are explored using maps and histograms, attached in the Appendix, specifically Supplemental Figures 13 to 30.

The eighteen rescaled variables are then converted from matrices sized 180x360 to eighteen strings of dimensions 1x64800. However, these strings are still riddled with NaNs from land pixels. To solve this problem, in a parallel process all the 2D matrices are stacked and added together, creating a comprehensive and gapless map of the latitude and longitude location of continental NaNs. This landmass index of NaNs is used to extract only the oceanic pixels from the strings, which result in eighteen strings of dimensions 1x41458. The landmass index of NaNs is put aside for later use. Finally, these strings are stacked vertically to create a matrix of dimensions 18x41458 (Figure 3). This is the input matrix used to run PCA.

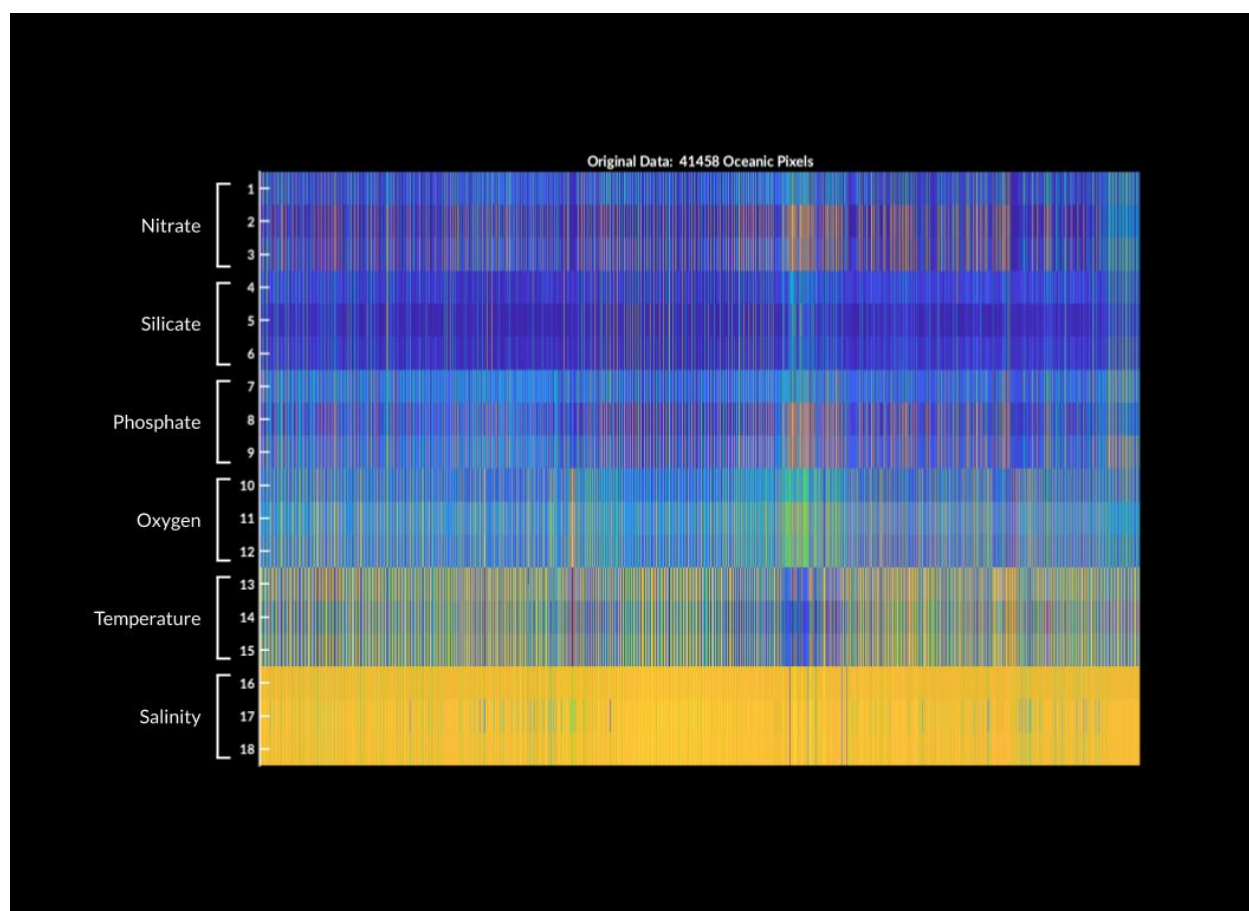


Figure 3. Input matrix of oceanic pixels in string form per oceanic factor for PCA.

1.3 CAPTURING VARIANCE

Running the 18x41458 input matrix through PCA results in three output matrices, the eigenvectors (18x18 matrix), the PCs (41458x18 matrix), and the eigenvalues (18x1 value column). The eigenvalues of these eighteen modes are converted to percent of variance explained and plotted. The first three modes explain 97% of total variability in the input matrix (Figure 4).

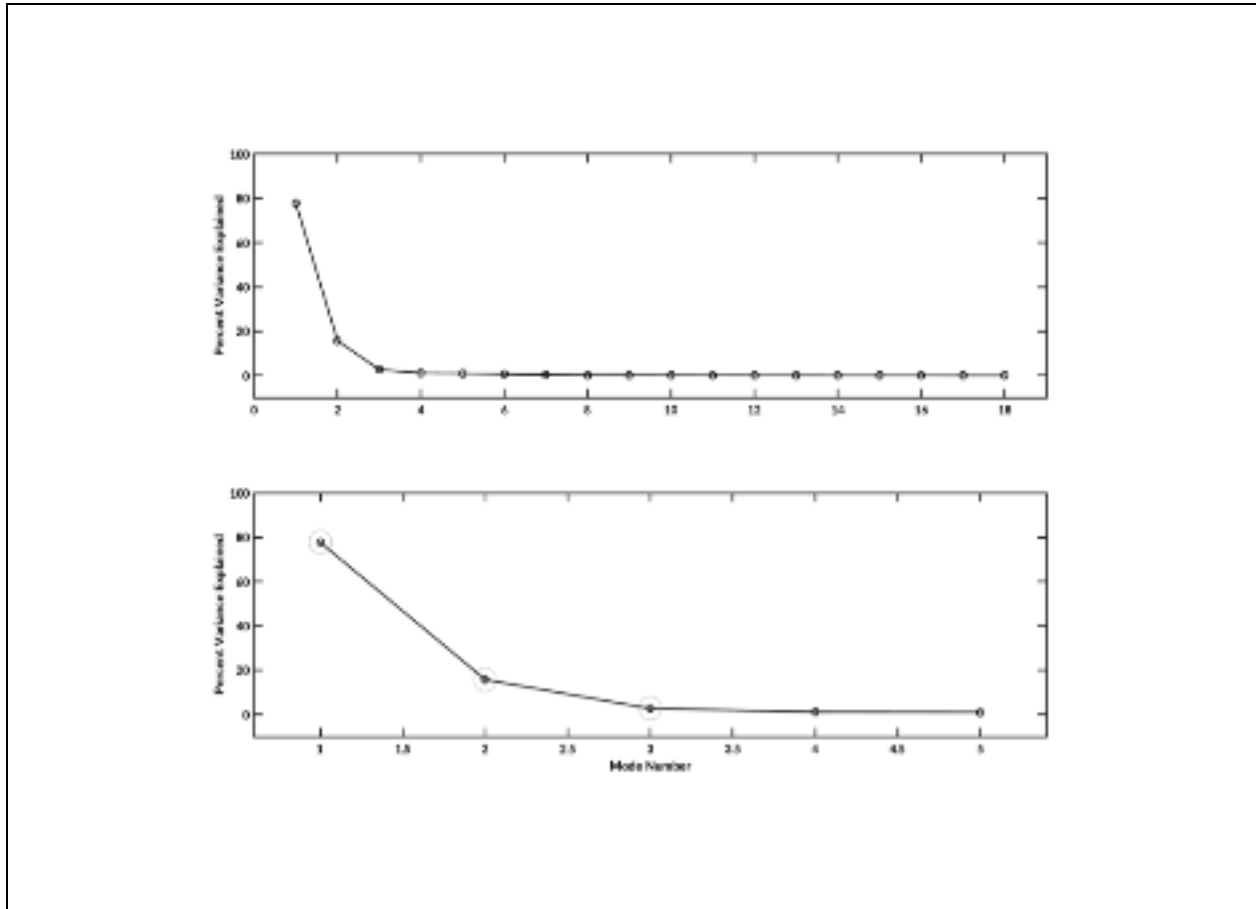


Figure 4. First three modes, where PC1 explains 78%, PC2 explains 16%, and PC3 explains 3%, for a total of 97% of the total variance explained.

The oceanic pixels of each PC can be converted from long string format into a 2D map matrix using the landmass index of NaNs, which was originally used to convert raw variable data into oceanic pixel strings. Because the first three modes capture close to 100% of variance, we will retain only these three for further analysis. Now we can look into what PC 1, PC 2, and PC 3 reveal about oceanic variance.

1.3.1 PC 1

The first mode explains 78% of total variance, and appears to be influenced by temperature and to a lesser extent by salinity (Figure 5). This is congruent with Earth's at-large oceanic physical processes, which are chiefly governed by the energy input of the Sun and by the local properties of seawater, like salinity and density. In other words, PC 1 is related to the energy stored in the physical aspect of the oceanic realm.

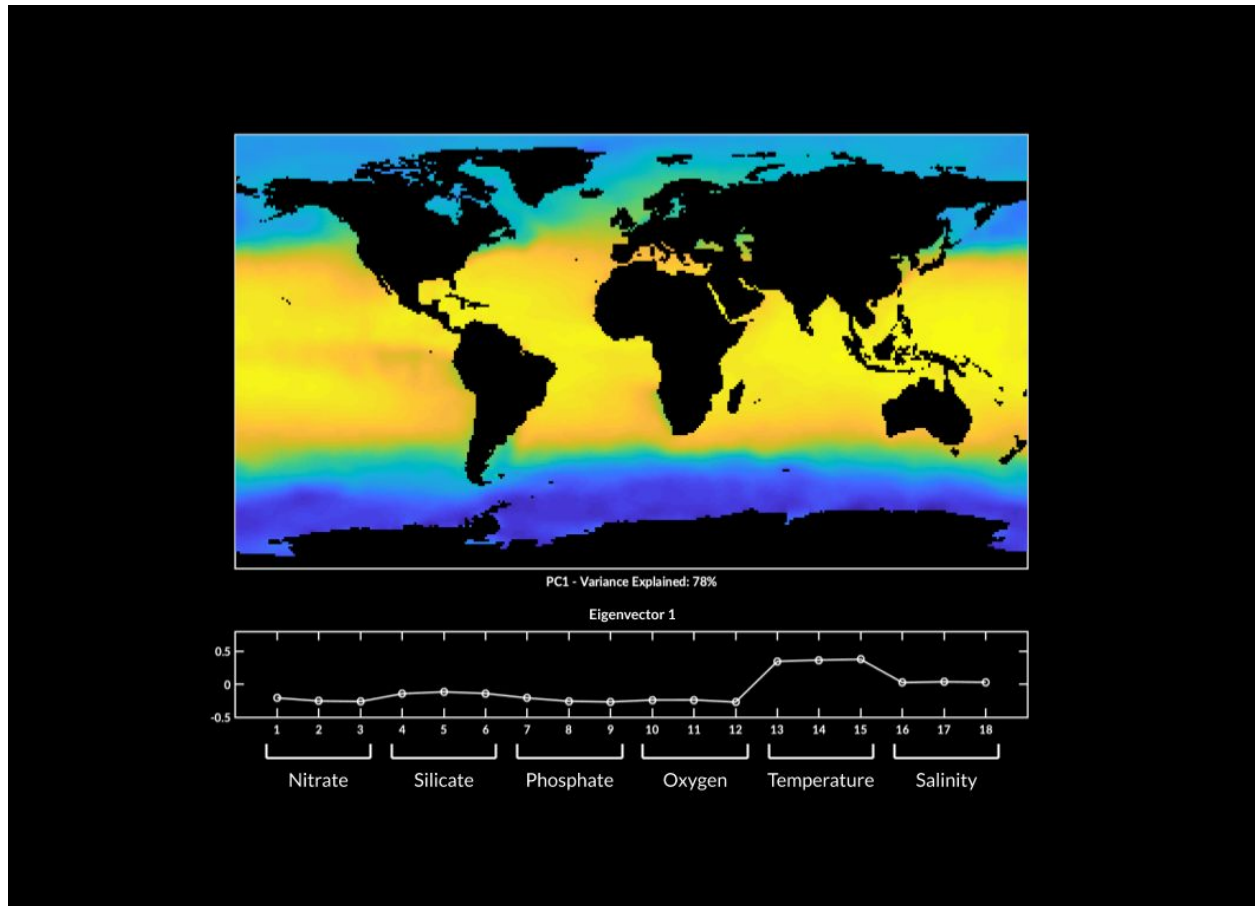


Figure 5. PCA Mode 1.

1.3.2 PC 2

The second mode explains 16% of the variance, and is influenced by several simultaneous signatures, but chiefly macronutrients (Figure 6). Granted temperature and salinity have been already accounted for in the first mode, the second mode appears to be influenced by macronutrient concentrations (nitrate, silicate, and phosphate), and inversely related to oxygen concentration.

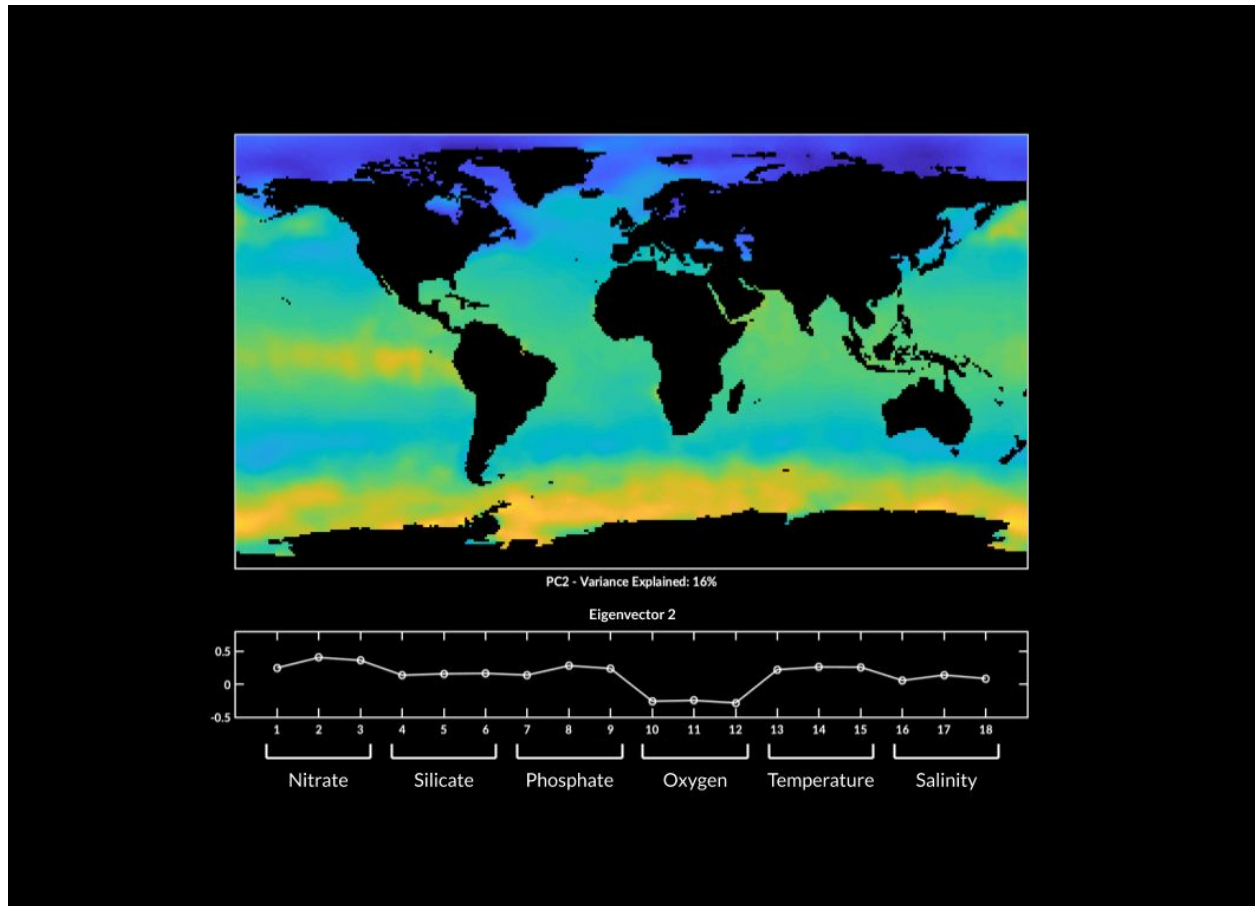


Figure 6. PCA Mode 2.

1.3.3 PC 3

The third mode explains 3% of the variance and is influenced by silicate (Figure 7). The strong mixing regions surrounding Antarctica upwell silicate to the oceanic surface where it can be uptaken by diatoms, which use silica for their glass shells. Because silicate is prominent in Antarctica, the geographic signature of this mode reveals important Southern Ocean features, such as the Polar Front.

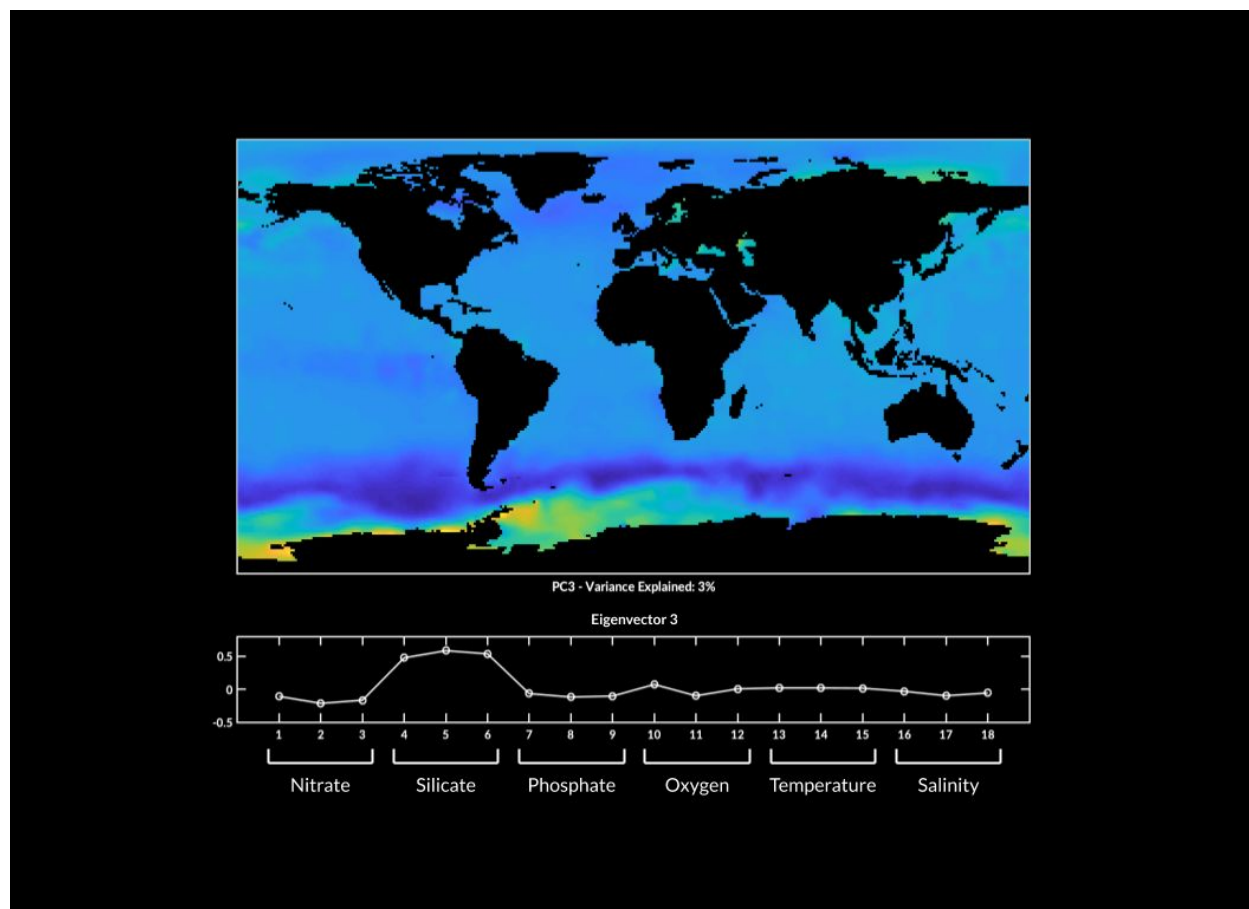


Figure 7. PCA Mode 3.

1.4 EXPLORING VARIANCE

The eigenvector matrix (18x18) and the PCs matrix (41458x18) can be plotted simultaneously to visualize how the eigenvectors explain the distribution of oceanic pixels in a particular combination of PCs. Because there are 6 factors subdivided into 3 components each (maximum, minimum, and mean), the vector themselves are averaged together so that a general factor vector is obtained. This results in only six vectors, one per factor, useful for a clearer visualization.

Figures 8, 9, and 10 show PC oceanic pixels and the six eigenvectors plotted side-by-side in duplicates of the same combination of PC axes, where Figure 8 shows the interaction of PC 1 and PC 2, Figure 9 shows the interaction of PC 1 and PC 3, and Figure 10 shows the interaction of PC 2 and PC 3.

Across the three combinations of PCs, there seem to be four general directions that govern distribution of oceanic pixels across principal component space. First are the macronutrients (nitrate, silicate, and phosphate), which align together and are expressed with greater intensity along the axis of PC 2. Silicate is separated from macronutrients, but only along the axis of PC 3. Second is oxygen, with points in the opposite direction to macronutrients along the axis of PC 2. However, on the axis of both PC 1 and PC 3 oxygen aligns with macronutrients. In contrast, third is temperature, which is orthogonal to all other eigenvectors, expressed with greatest intensity along the axis of PC 1. Finally, fourth, is salinity, which does not have prevalence across any PC axis, so it may be regarded as minimally influential in its contribution to overall variance.

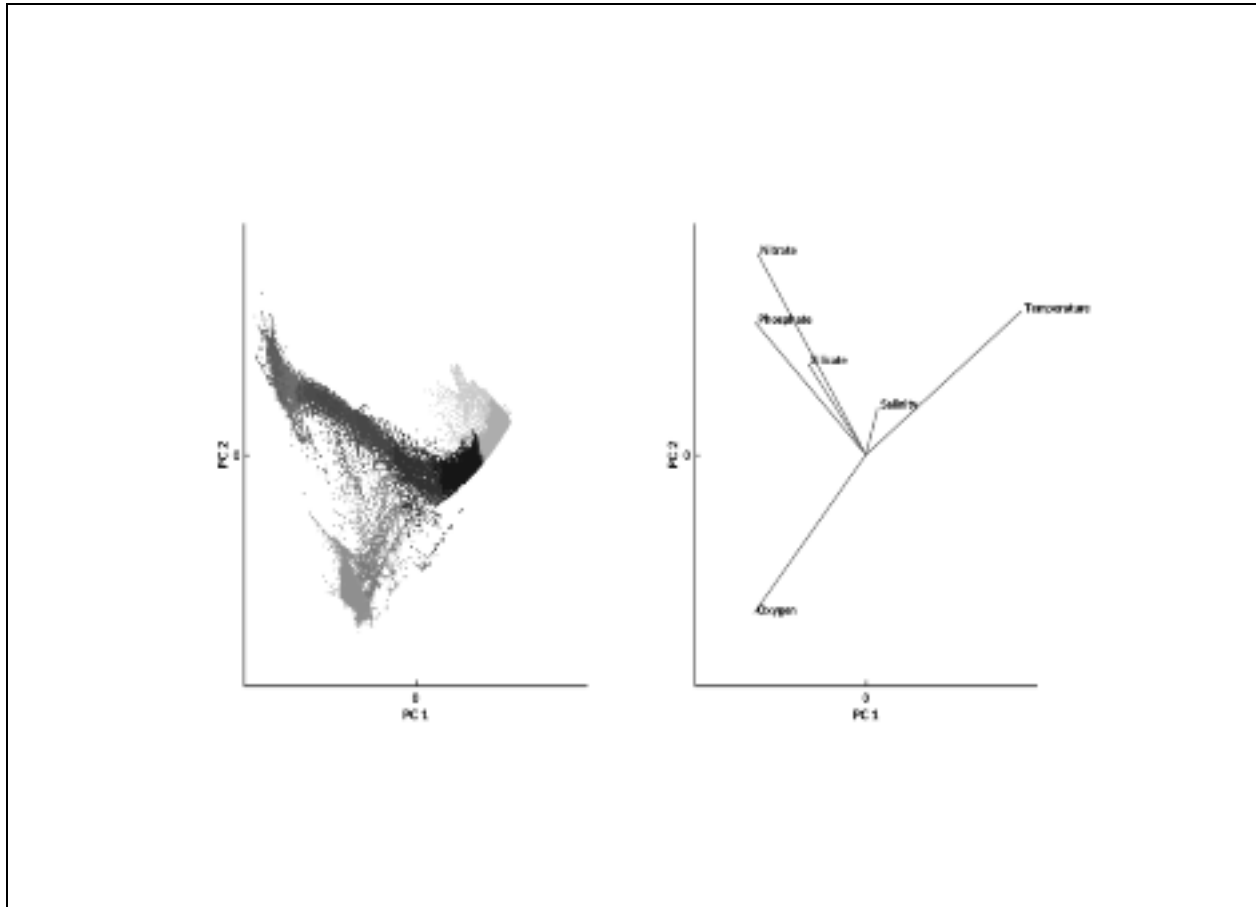


Figure 8. Oceanic pixels and variable vectors plotted PC 1 and PC 2 space.

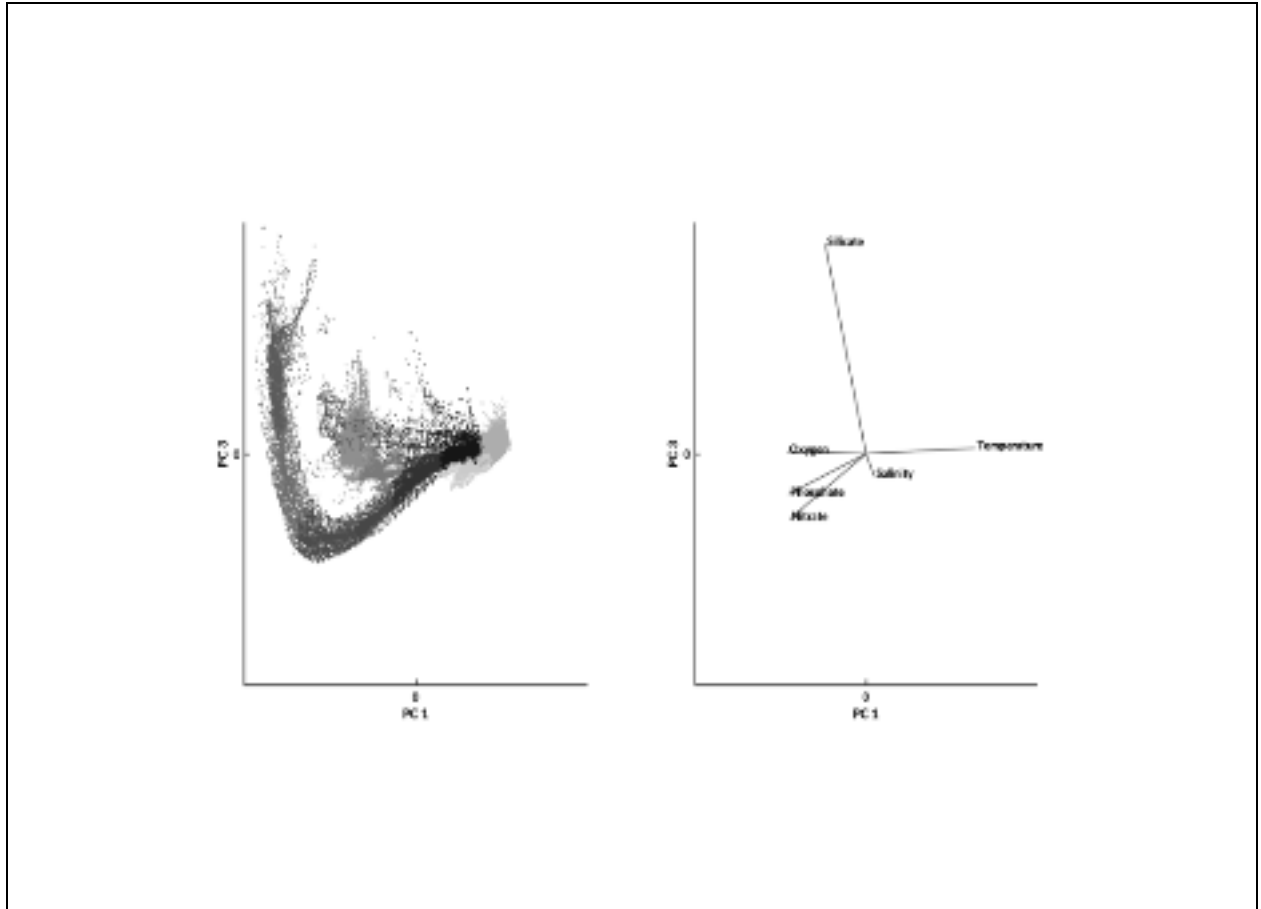


Figure 9. Oceanic pixels and variable vectors plotted PC 1 and PC 3 space.

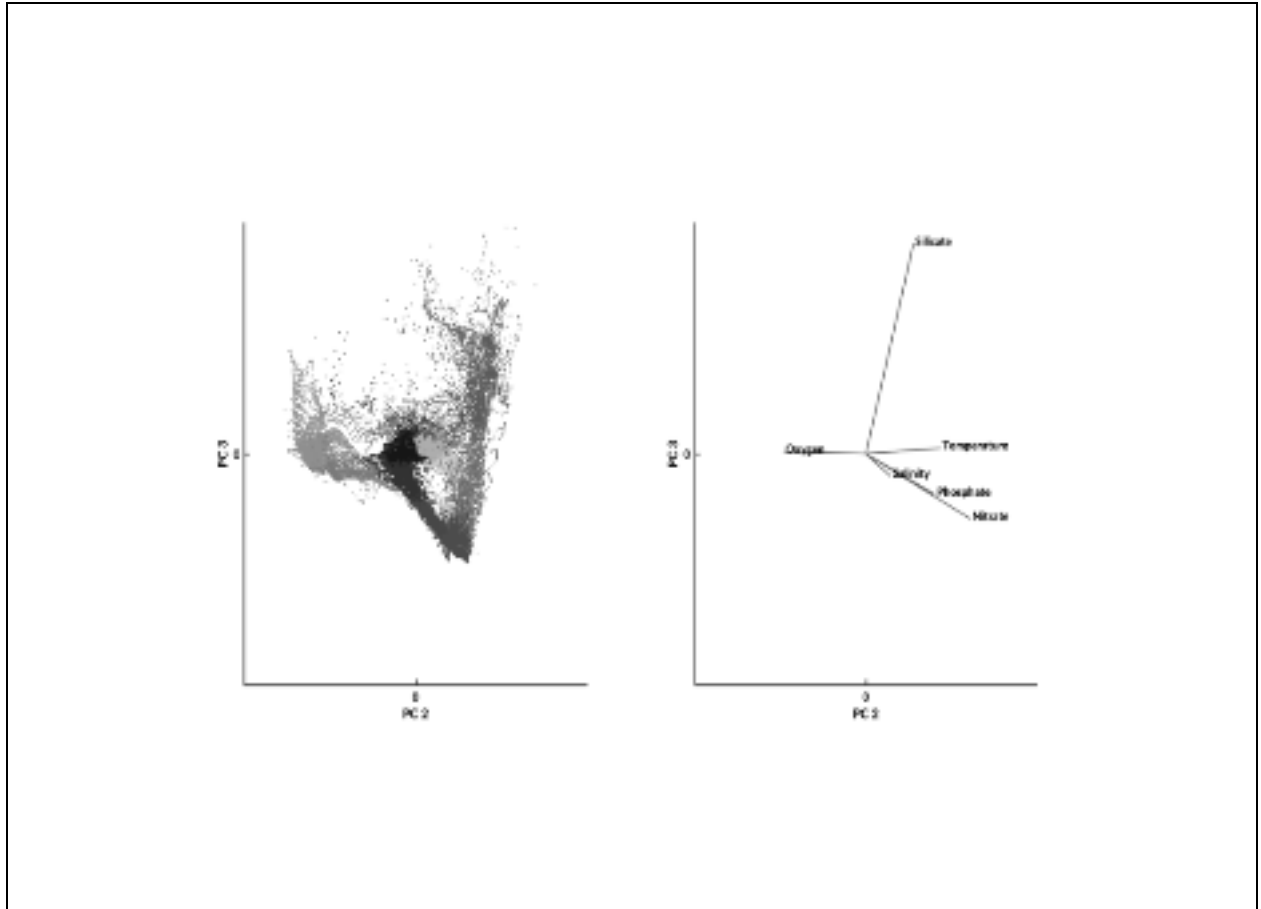


Figure 10. Oceanic pixels and variable vectors plotted PC 2 and PC 3 space.

1.5 HIERARCHICAL CLUSTERING

After exploring the results of PCA, the first three modes are now hierarchically clustered using Ward's method with euclidean distance, and the results are plotted in a dendrogram (Figure 11). Based on the dendrogram, five clusters or $k=5$ is a good place to start looking at what patterns clustering reveals.

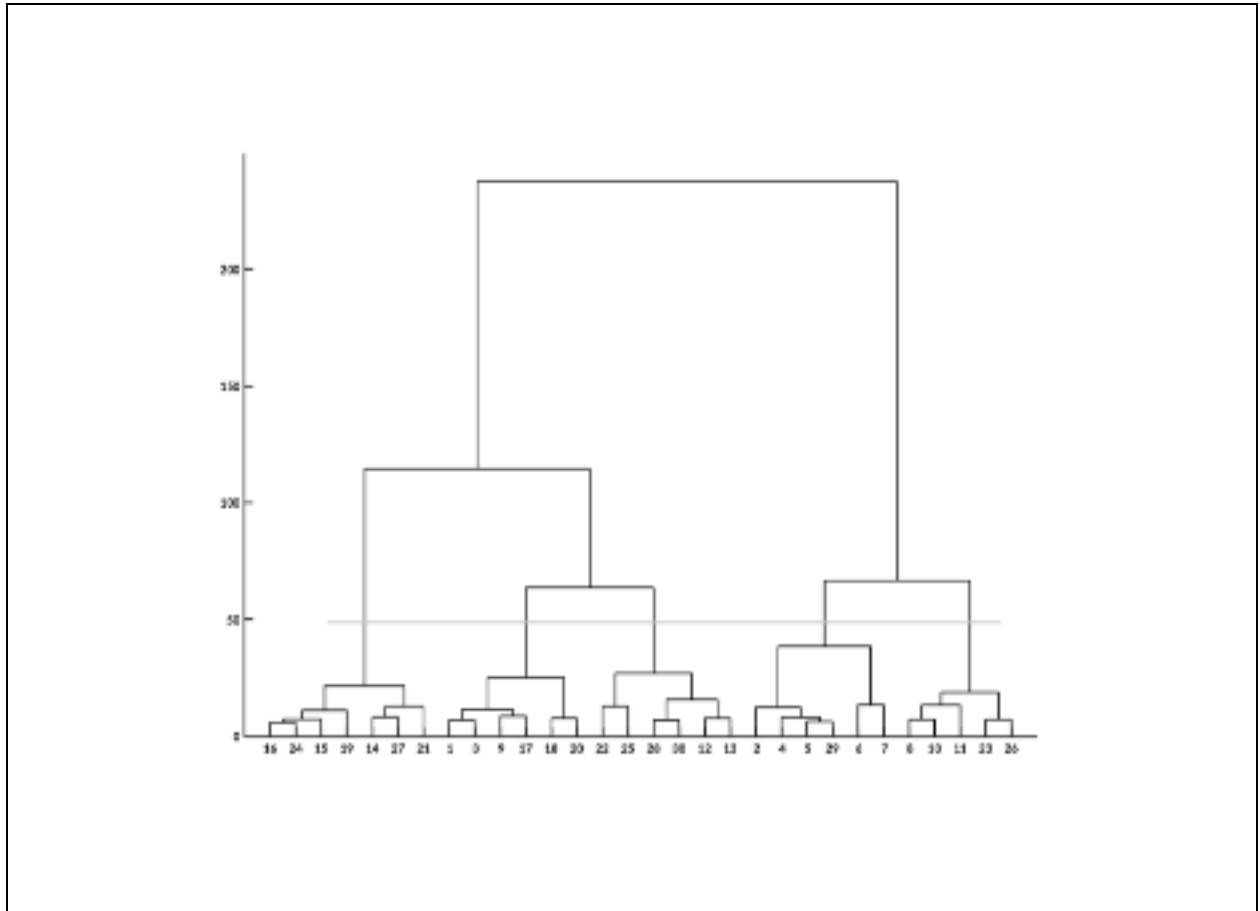


Figure 11. Linkage dendrogram from hierarchical clustering.
The horizontal line represents the threshold for 5 clusters, or $k=5$.

2. OCEANIC ECOREGIONS

2.1 WORLD BANDS

Starting analysis at $k=5$ yields 5 clusters, which would correspond to oceanic ecoregions at the whole-planet scale, better understood as biomes. Because of its tilt, if Earth was purely oceanic the surface would have planetary bands around the tropics and poles (Parsons and Lalli, 2006). At 5 clusters these ecologically-meaningful bands appear, which correspond to the major oceanic biomes determined by latitude. These are Tropical Waters, Subtropical Waters, Subpolar Waters, Antarctic Waters, and Arctic Waters.

Interestingly, there is a pocket of Antarctic-like water in the North Pacific that has no equivalent in the North Atlantic. This may be because the Atlantic is a younger ocean, whereas the Pacific is still the original Panthalassic Ocean, but this is just speculation. Furthermore, the Arctic ecoregion is identified early in the clustering process, and persists unrelated to other ecoregions and robust against further clustering.

There are many ways to represent oceanic clusters. First, Figure 12 shows a map is composed of oceanic pixels reconstructed using the landmass index of NaNs. This map intuitively shows the world band ecoregions that arise from $k=5$. These clusters are also visible in PC space, as shown in Figure 13. However, a shortcoming of classical map projections, like the Lambert-like projection shown in Figure 12, is that it is landmass-centric, interrupting the ocean. To address this, Figure 14 shows the same world bands as they appear across the entire World Ocean, uninterrupted.

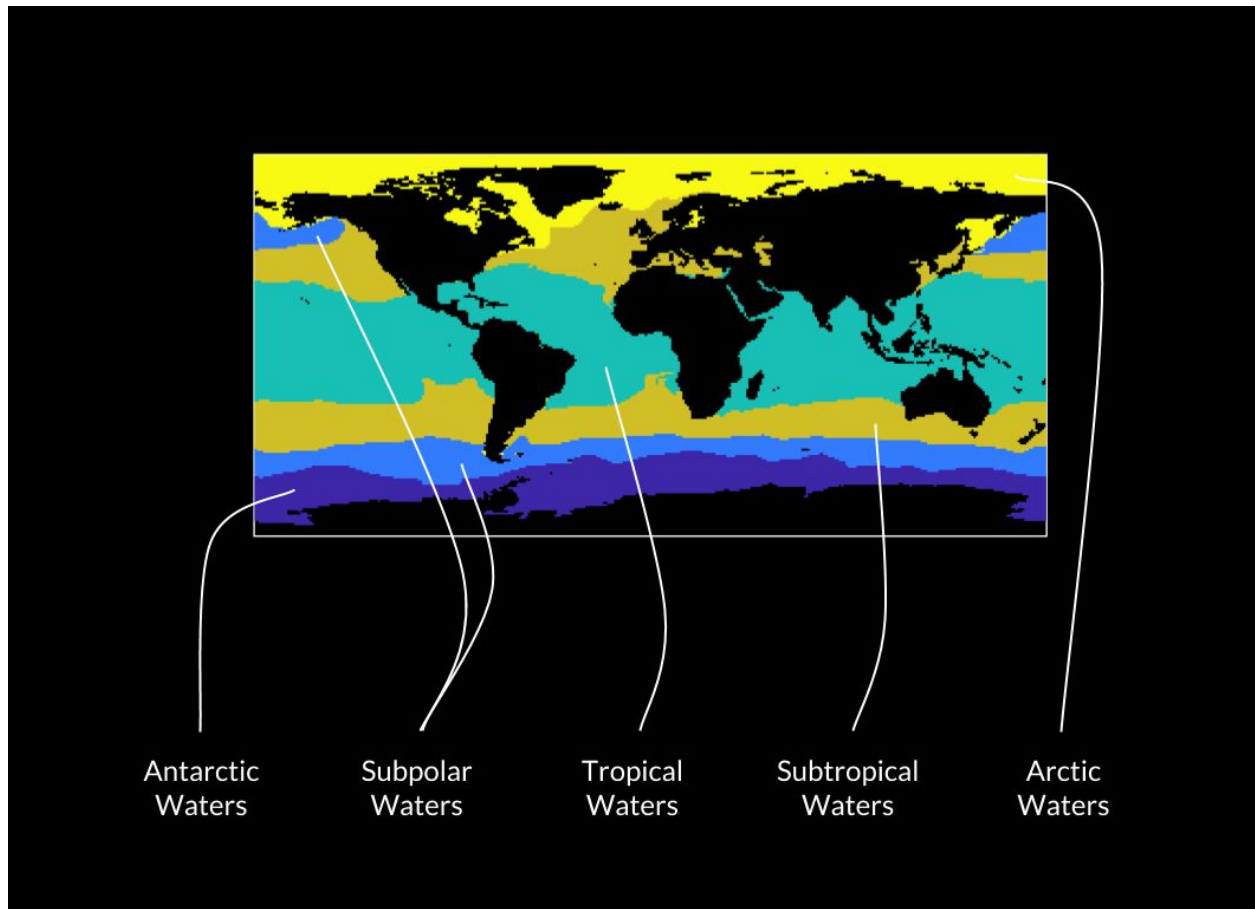


Figure 12. Five oceanic biomes, in this case world bands, identified by $k=5$.
The map is composed of oceanic pixels reconstructed as from a string of pixels that results from hierarchical clustering.

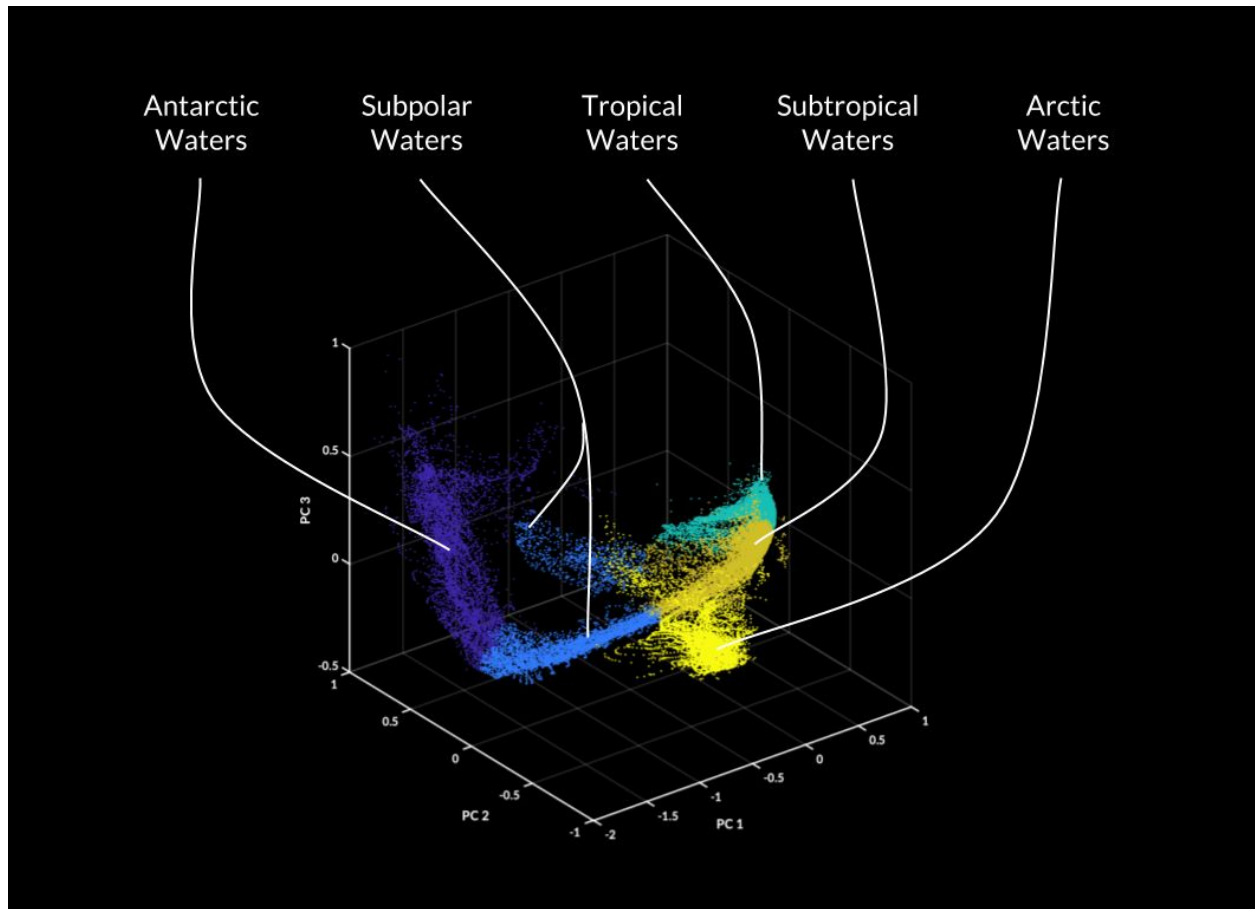


Figure 13. Five oceanic biomes, in this case world bands, identified by $k=5$. Here oceanic pixels from hierarchical clustering are plotted in the 3D principal component space of the first three PCs.

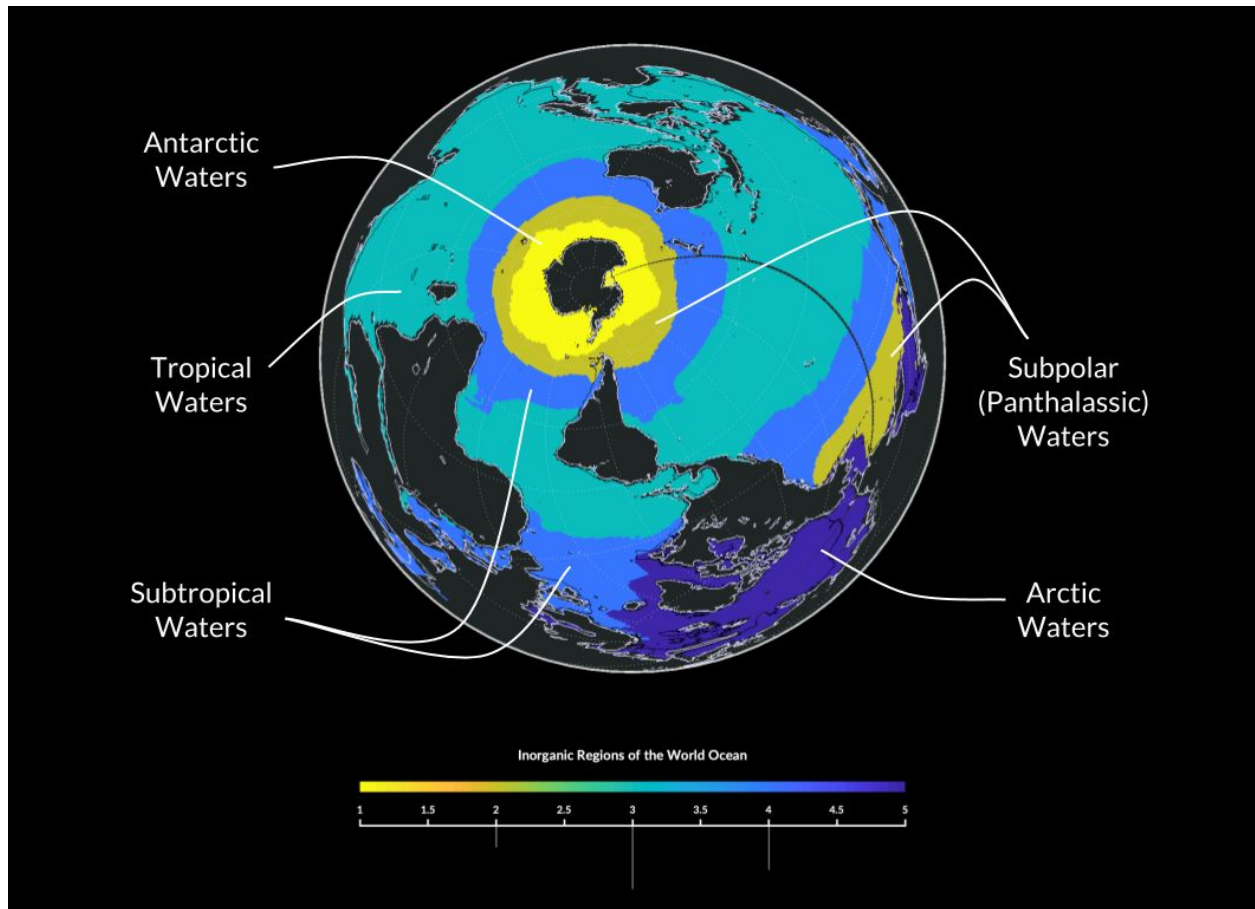


Figure 14. Five oceanic biomes, in this case world bands, identified by $k=5$. The map is composed of oceanic pixels reconstructed as from a string of pixels that results from hierarchical clustering, plotted in an azimuthal equidistant map where the center of projection is the antipode of the landlocked pole of inaccessibility within Eurasia, which results in the World Ocean shown completely and uninterrupted.

2.2 CLUSTERS WITHIN CLUSTERS

By increasing the threshold k , ecoregions systematically dissociate into smaller ecoregions. At 10 clusters, or $k=10$, the great tropical band (Figure 14, Tropical Waters) finally breaks and reveals western seafront continental upwelling regions, famous for their high productivity. Upwelling regions are macronutrient rich, so they produce blooms of phytoplankton that then precipitate as marine snow from the sunlit surface. Marine snow is then consumed alongside oxygen at depth. Hypoxic remineralization of iron into pyrite releases oxygen, so there is a surface signature over the oxygen minimum zone, or OMZ, in the shape of the hypoxic region (Parsons and Lalli, 2006). The Equatorial Pacific, Benguela Current, and Arabian Sea are known regions of hypoxic water, and are revealed at 10 clusters (Figure 15).

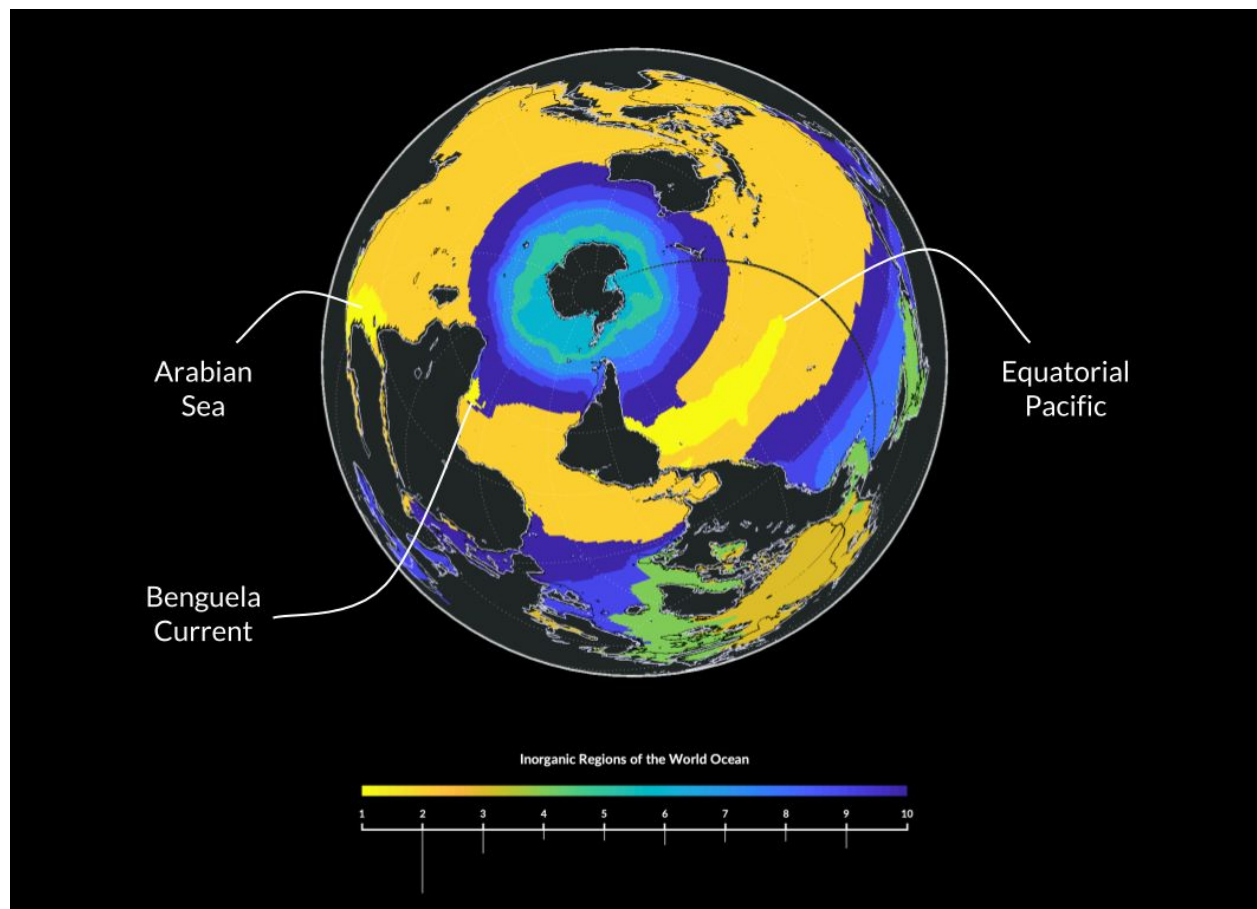


Figure 15. Ecoregions that emerge at $k=10$.

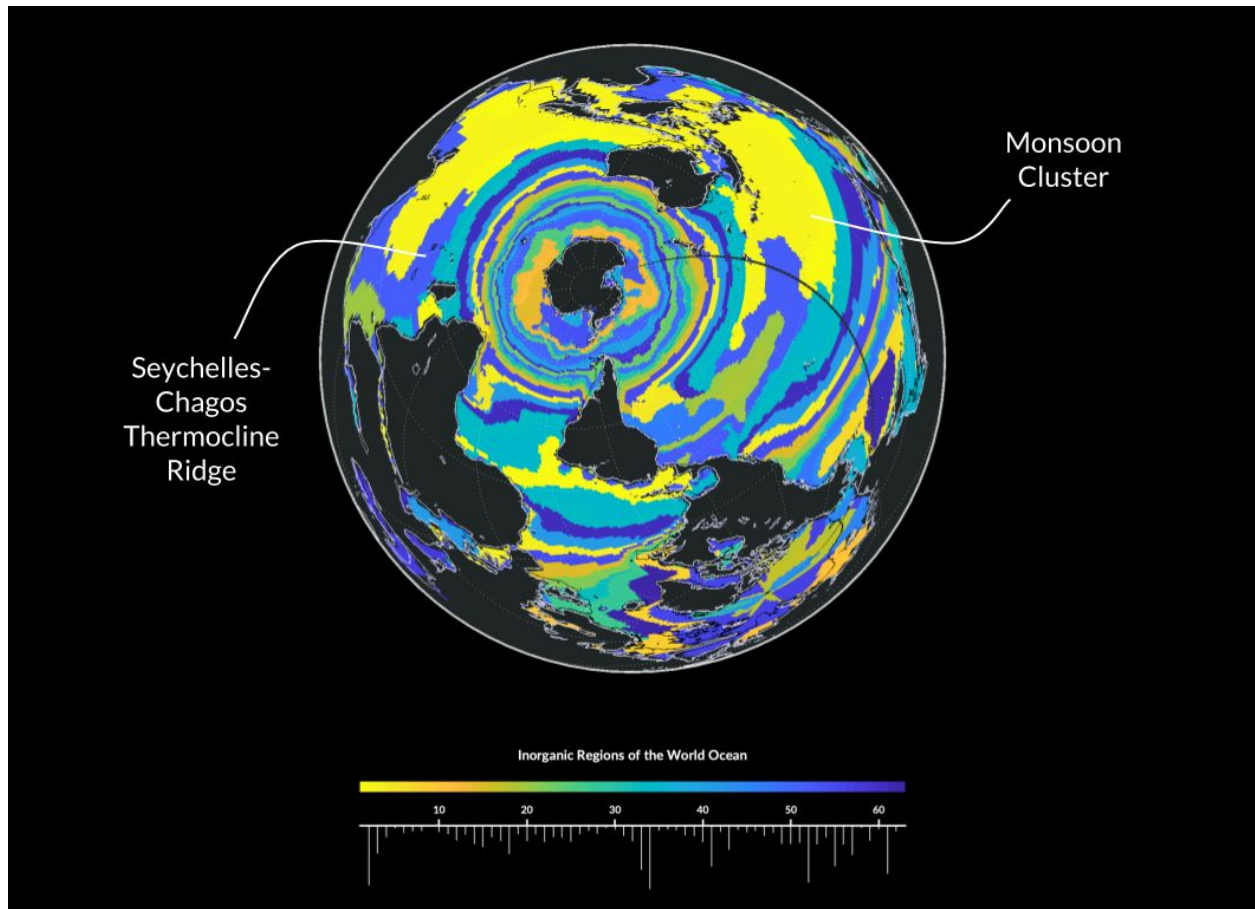


Figure 16. Ecoregions that emerge at $k=63$.

Some regions are persistent, for example, even when clustered beyond $k=60$, a very large cluster that includes the waters south of India, Southeastern Asia, and Northeastern Oceania remains unbroken. Because of its importance in climate, this oceanic region could be called the Monsoon Cluster. Only at 63 clusters this region finally breaks, revealing subdivisions around the edges (Figure 16). One of these regions resembles the open-ocean upwelling region known as the Seychelles-Chagos thermocline dome.

2.3 BOUNDARY STRENGTH

For every new k value a new cluster breaks apart from the collection of existing clusters. This newest cluster is the contribution of a given k value to overall partitioning. By systematically stacking the boundary of this new cluster, it is possible to determine which oceanic pixels are part of strong ecoregion boundaries, meaning they belong to the edge of an ecoregion. Conversely, pixels that experience no update from k iterations are robust, meaning they belong to the core of an ecoregion. Figure 17 illustrates the systematic update of the newest cluster boundary per value of k . The results of stacking these boundaries is shown in Figure 18.

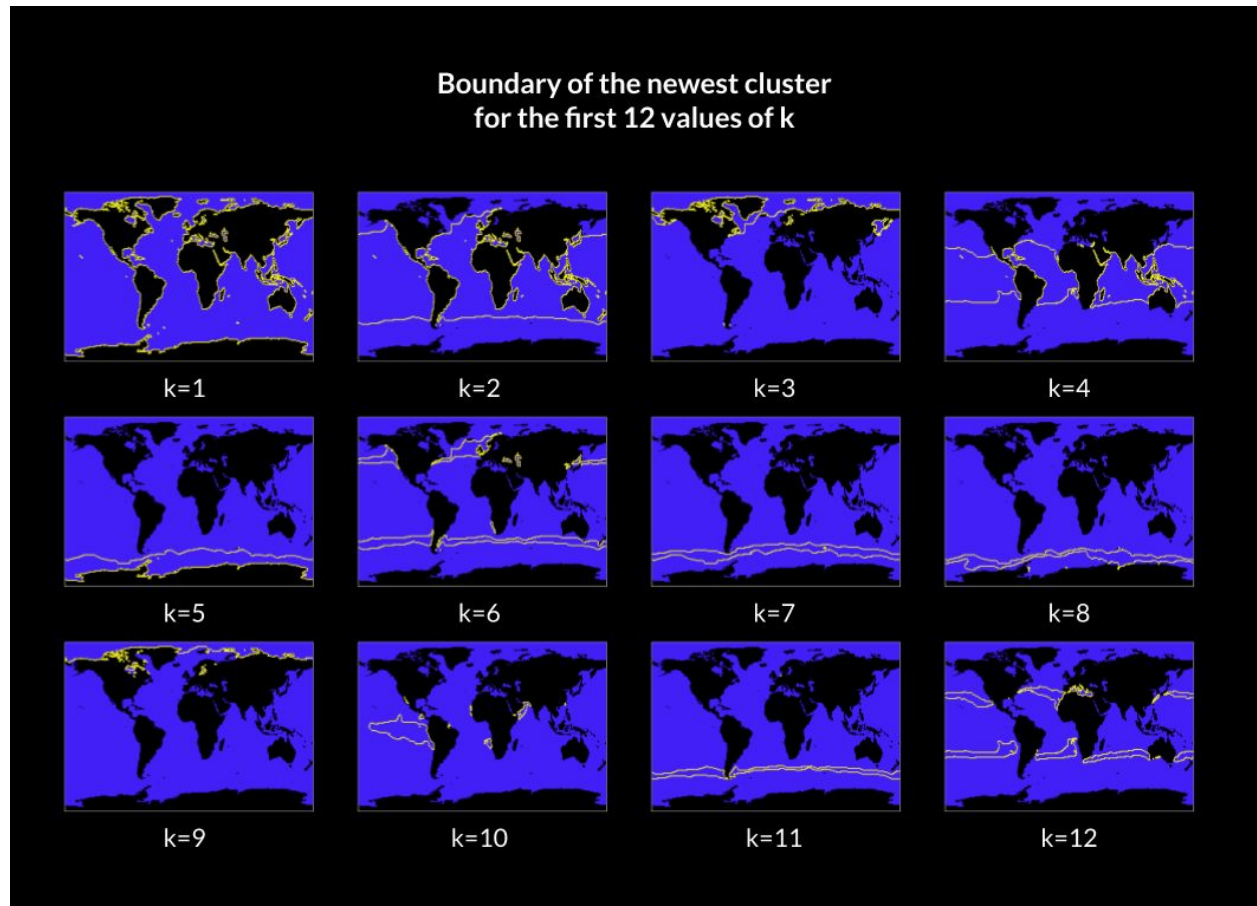


Figure 17. Newest boundary for each iteration of k , here shown are the first twelve values of k .

Strength of ecoregion boundaries at 2000 clusters

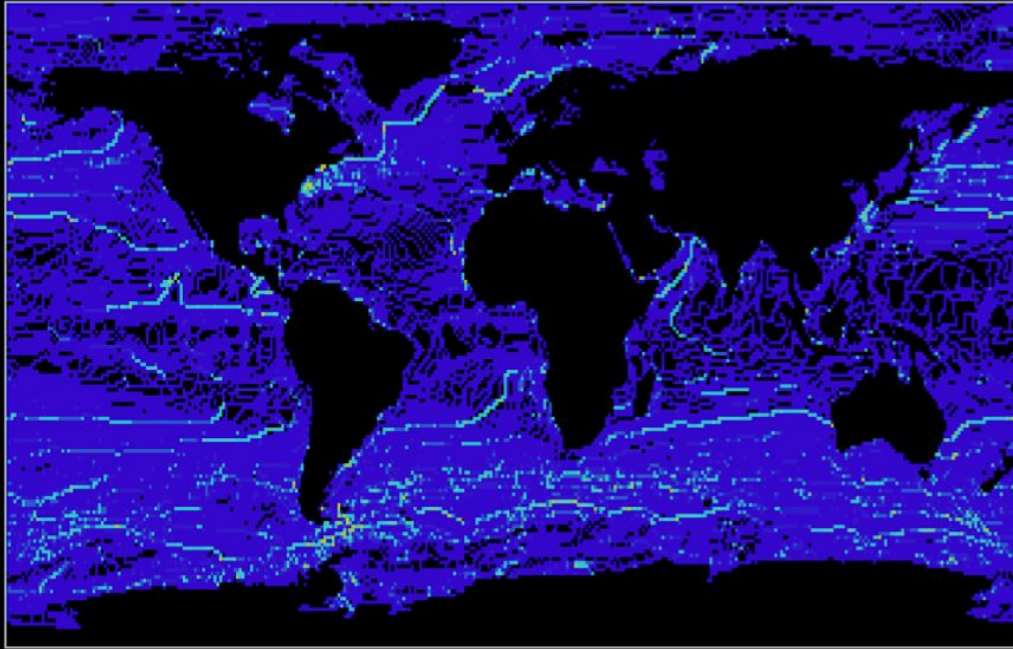


Figure 18. Ecoregion boundary strength detected after sequentially stacking boundary pixels from $k=1$ to $k=2000$.

Interestingly, at 2000 stacked clusters, the boundaries two elusive oceanic ecoregions appear, the Pacific Ocean Costa Rica thermocline dome and the Indian Ocean Seychelles-Chagos thermocline ridge (Figure 19).

Strength of ecoregion boundaries at 2000 clusters

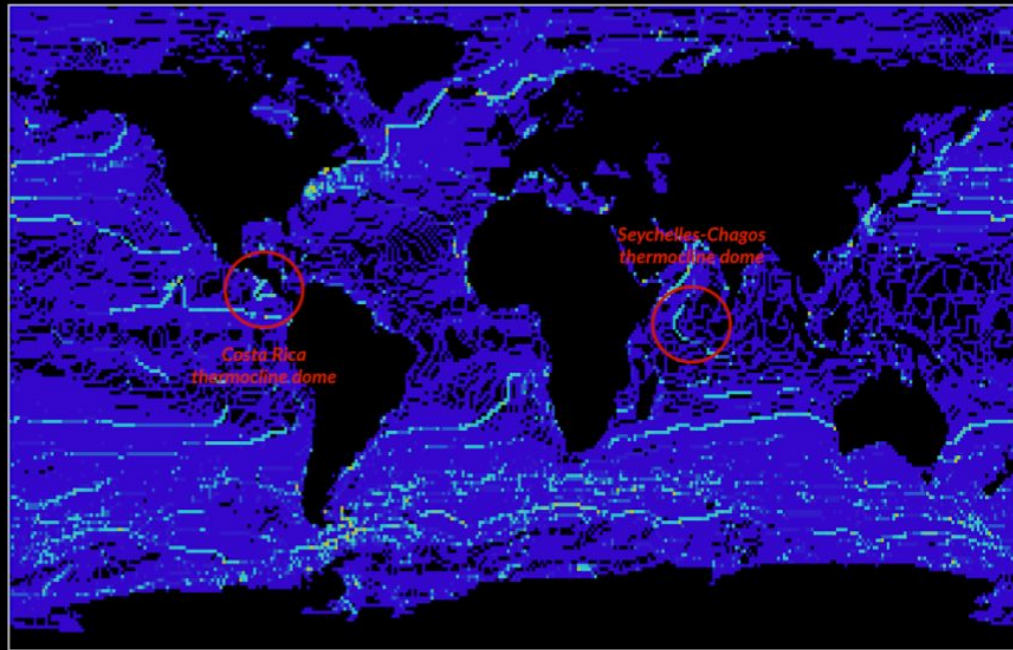
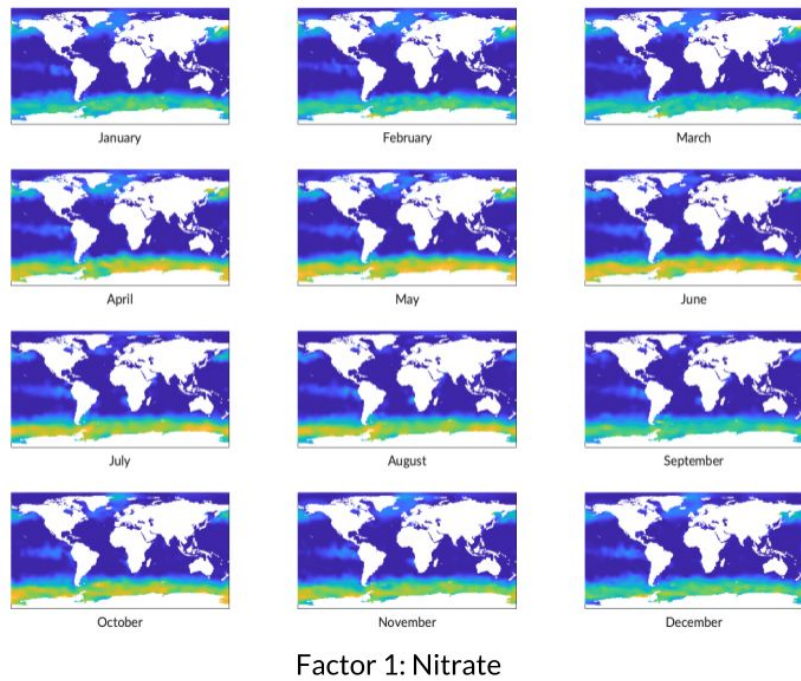


Figure 19. Thermocline dome boundaries.

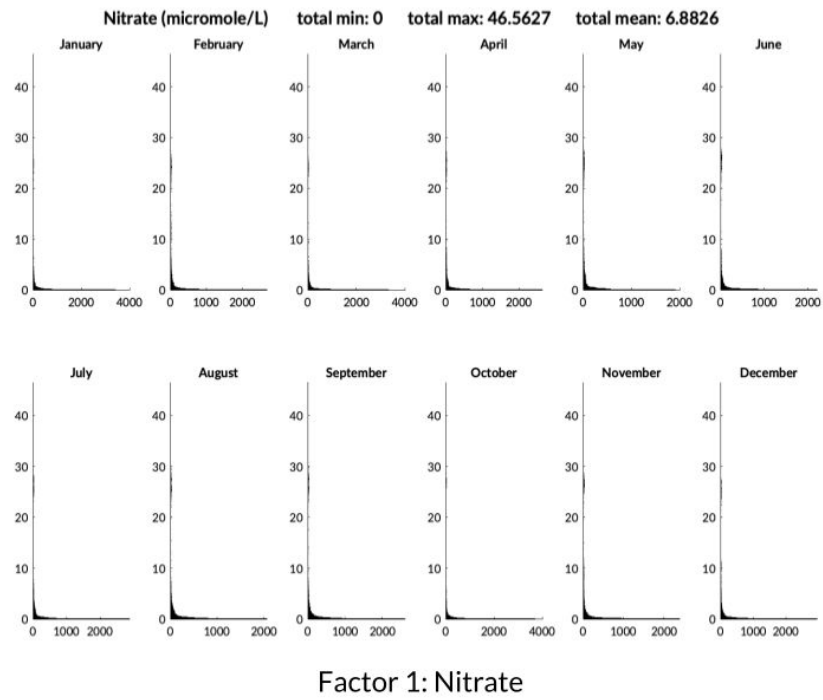
5. REFERENCES

Parsons and Lalli, 2006. Biological oceanography: An introduction.
WOA2013: <https://www.nodc.noaa.gov/OC5/woa13/woa13data.html>

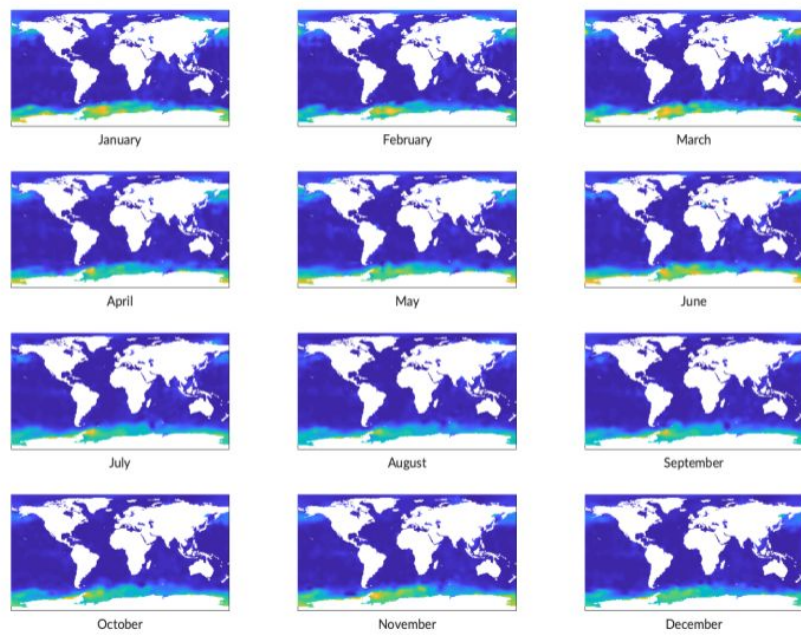
APPENDIX



Supplemental Figure 1. Maps of spatio-temporal variability for Nitrate (Factor 1).

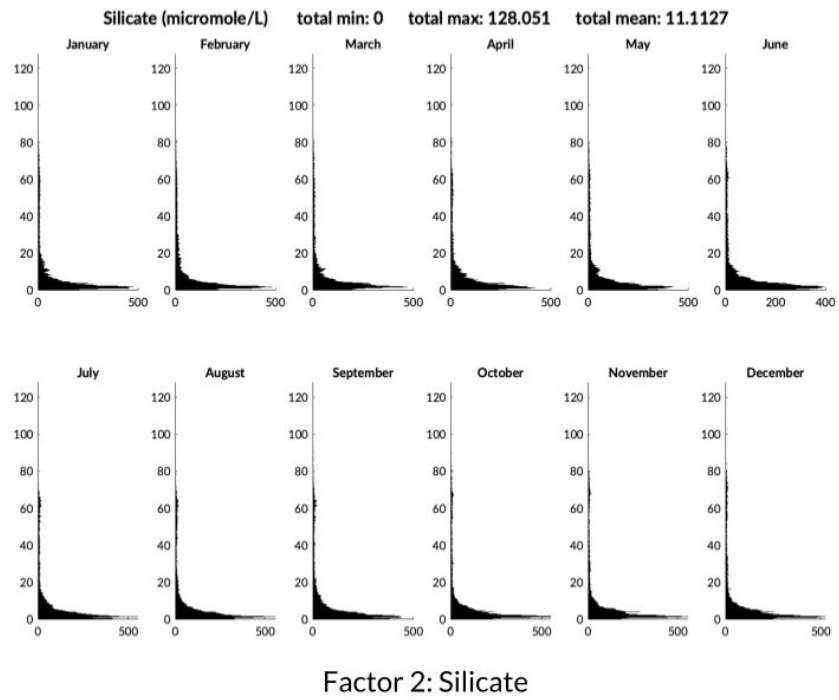


Supplemental Figure 2. Histogram of spatio-temporal variability for Nitrate (Factor 1).

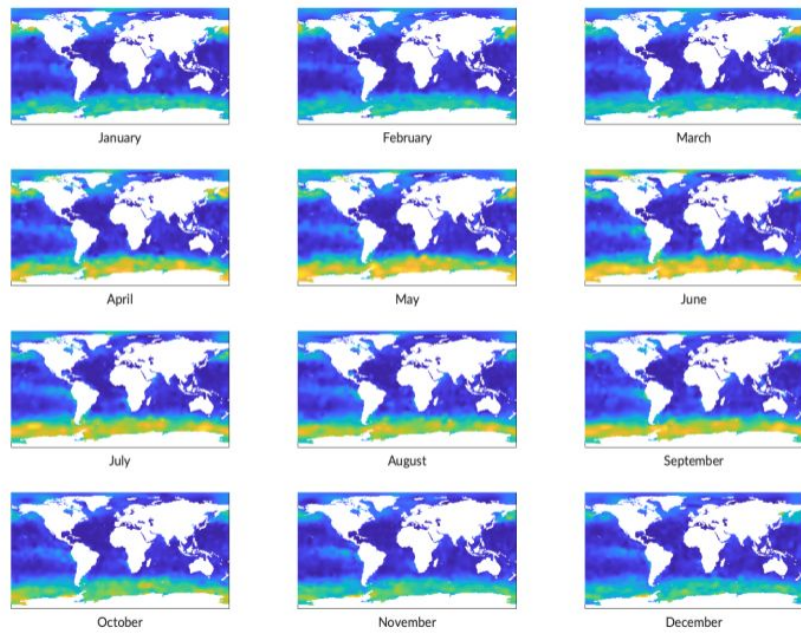


Factor 2: Silicate

Supplemental Figure 3. Maps of spatio-temporal variability for Silicate (Factor 2).

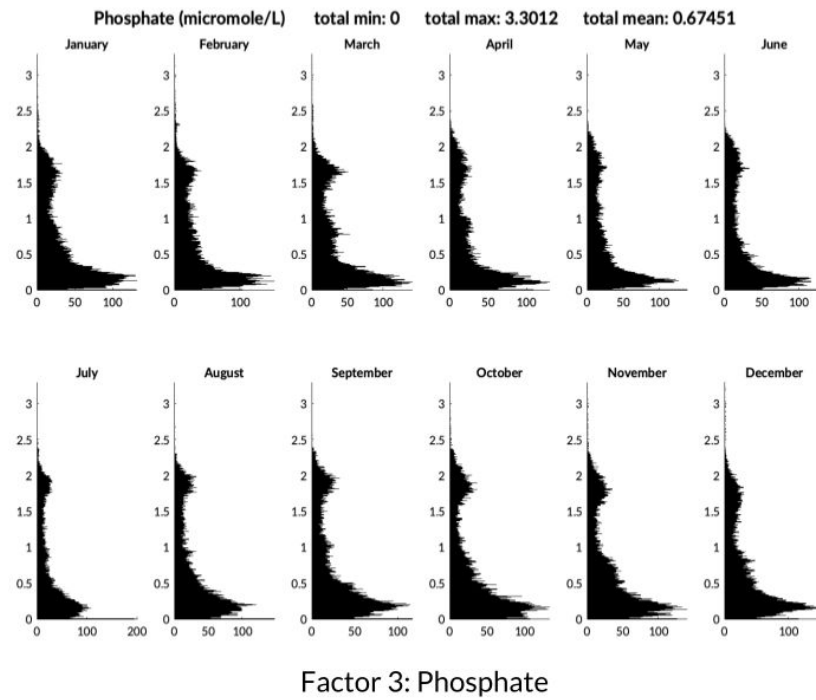


Supplemental Figure 4. Histogram of spatio-temporal variability for Salinity (Factor 2).

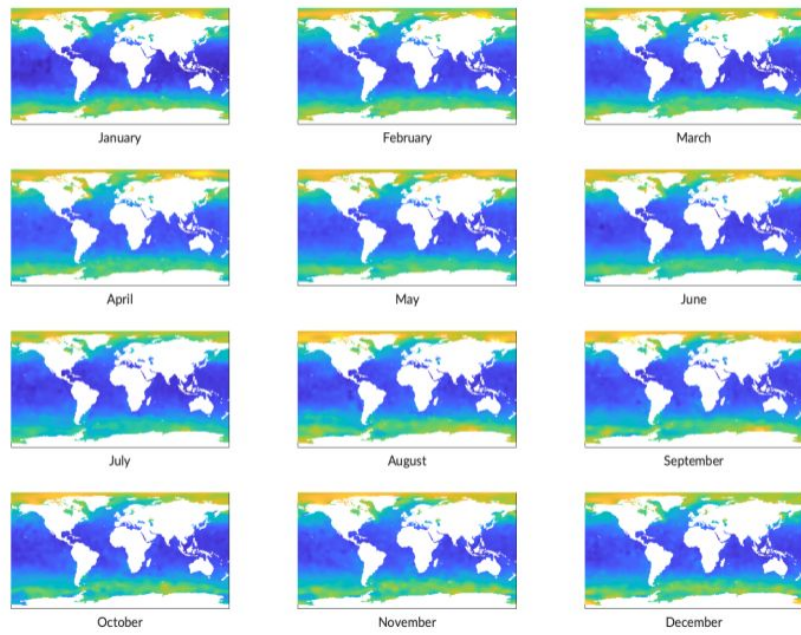


Factor 3: Phosphate

Supplemental Figure 5. Maps of spatio-temporal variability for Phosphate (Factor 3).

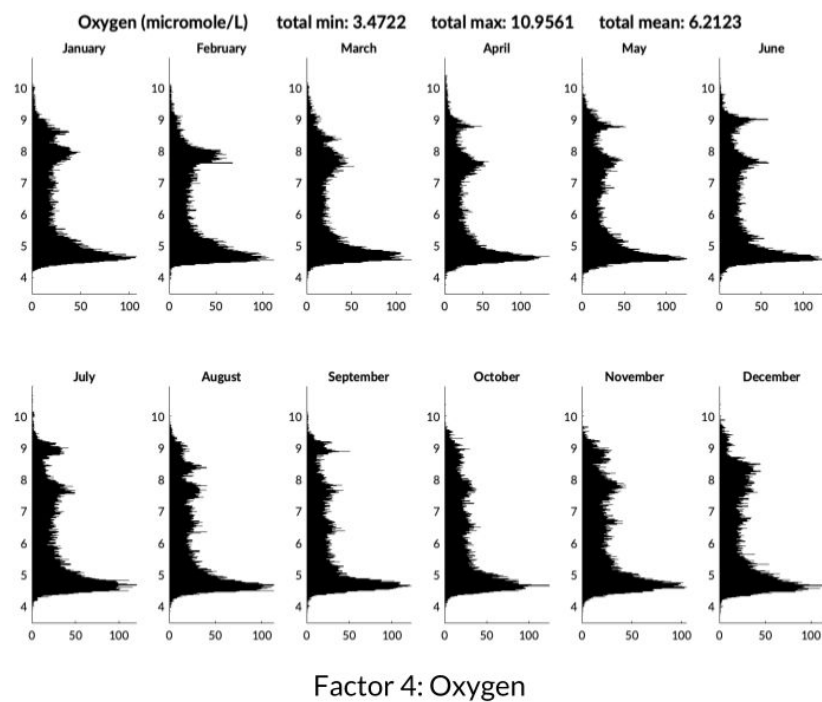


Supplemental Figure 6. Histogram of spatio-temporal variability for Phosphate (Factor 3).

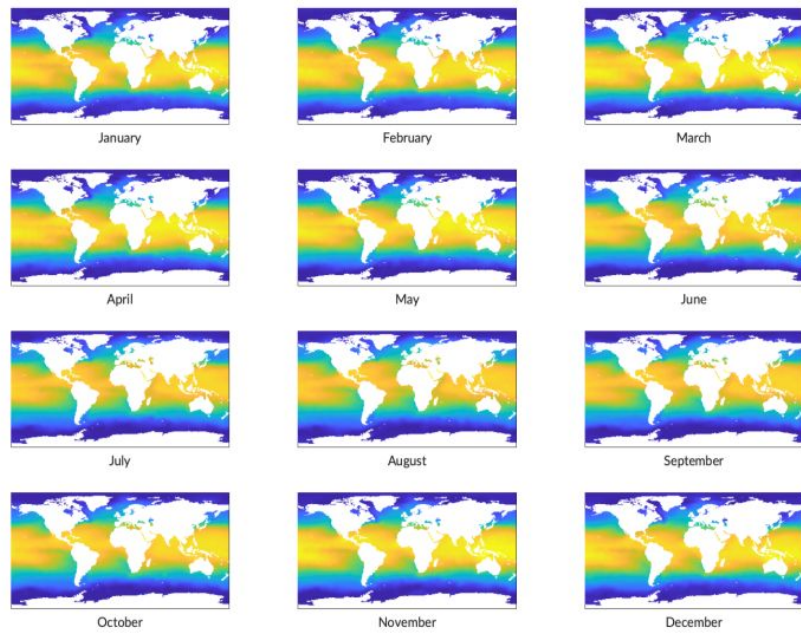


Factor 4: Oxygen

Supplemental Figure 7. Maps of spatio-temporal variability for Oxygen (Factor 4).

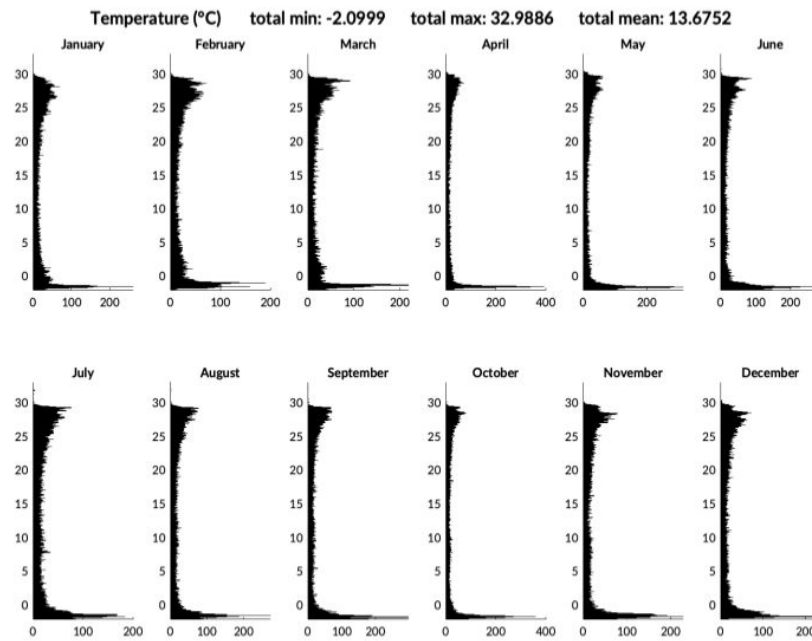


Supplemental Figure 8. Histogram of spatio-temporal variability for Oxygen (Factor 4).



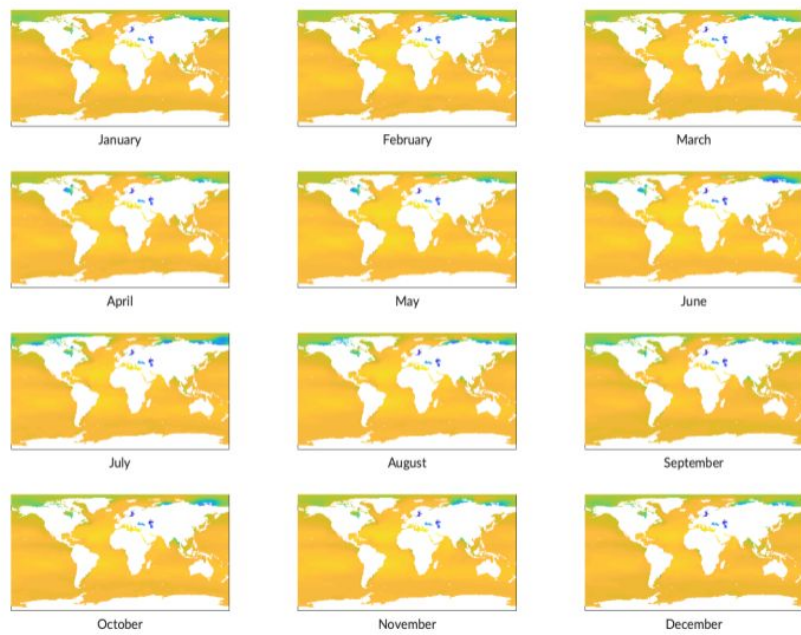
Factor 5: Temperature

Supplemental Figure 9. Maps of spatio-temporal variability for Temperature (Factor 5).



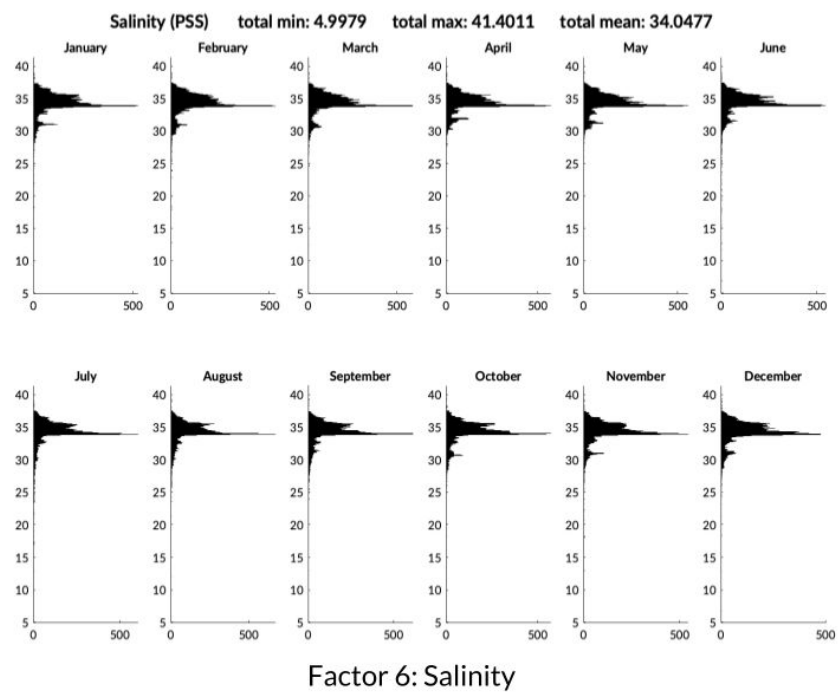
Factor 5: Temperature

Supplemental Figure 10. Histogram of spatio-temporal variability for Temperature (Factor 5).

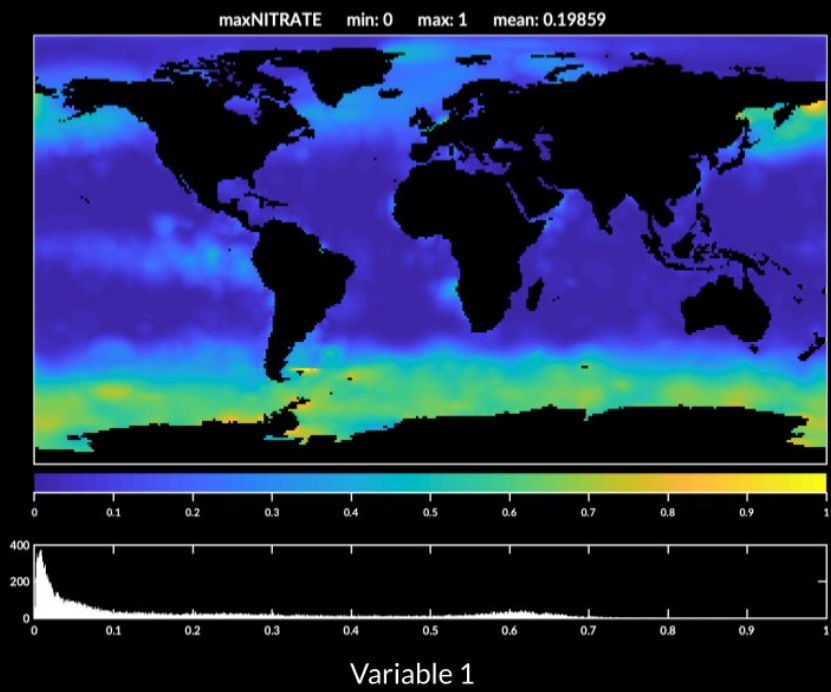


Factor 6: Salinity

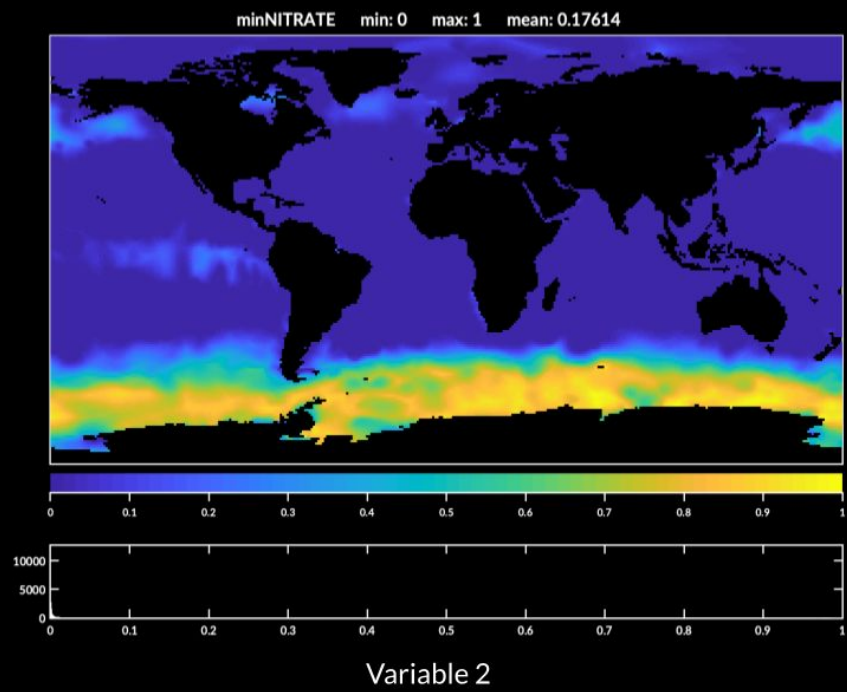
Supplemental Figure 11. Maps of spatio-temporal variability for Salinity (Factor 6).



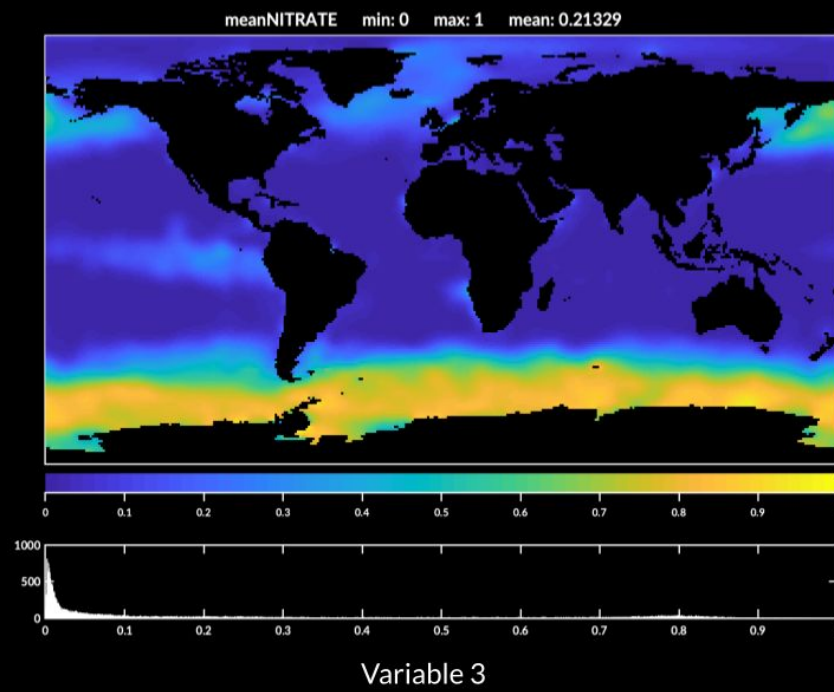
Supplemental Figure 12. Histogram of spatio-temporal variability for Salinity (Factor 6).



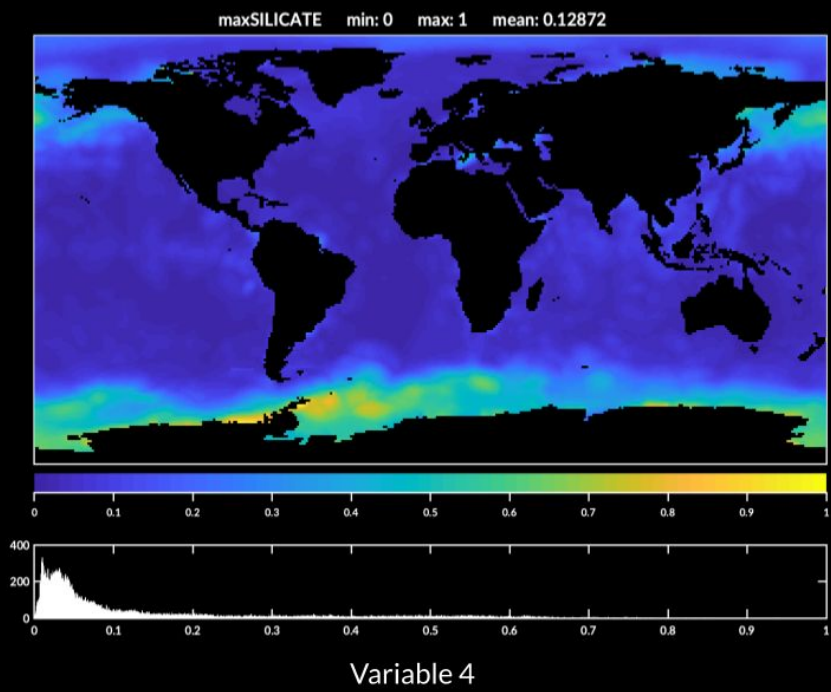
Supplemental Figure 13. Variable 1.



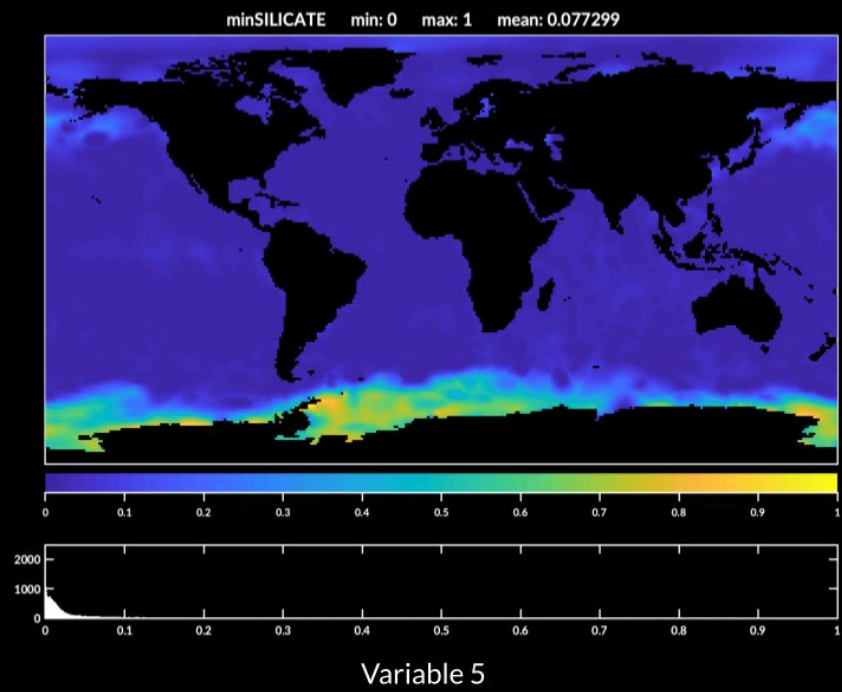
Supplemental Figure 14. Variable 2.



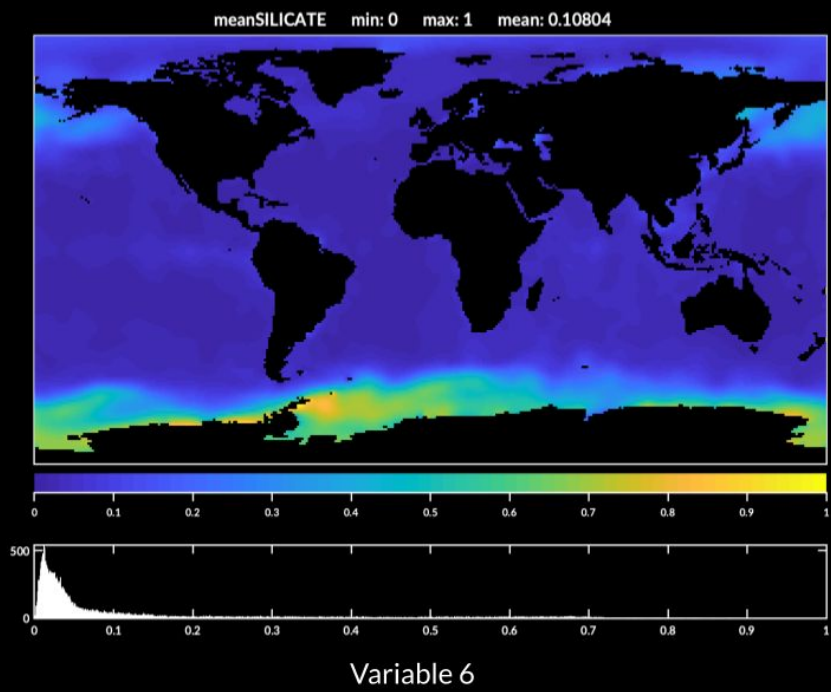
Supplemental Figure 15. Variable 3.



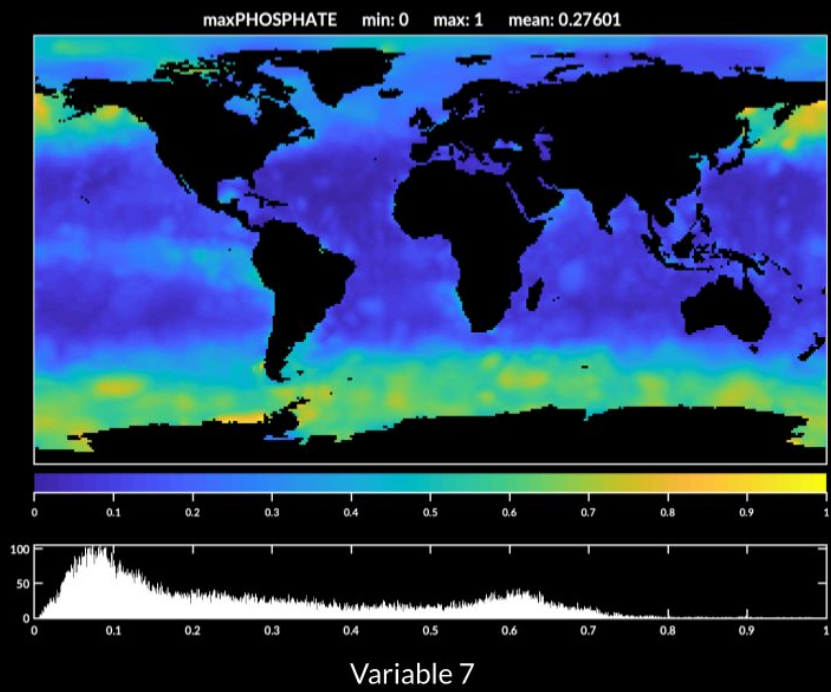
Supplemental Figure 16. Variable 4.



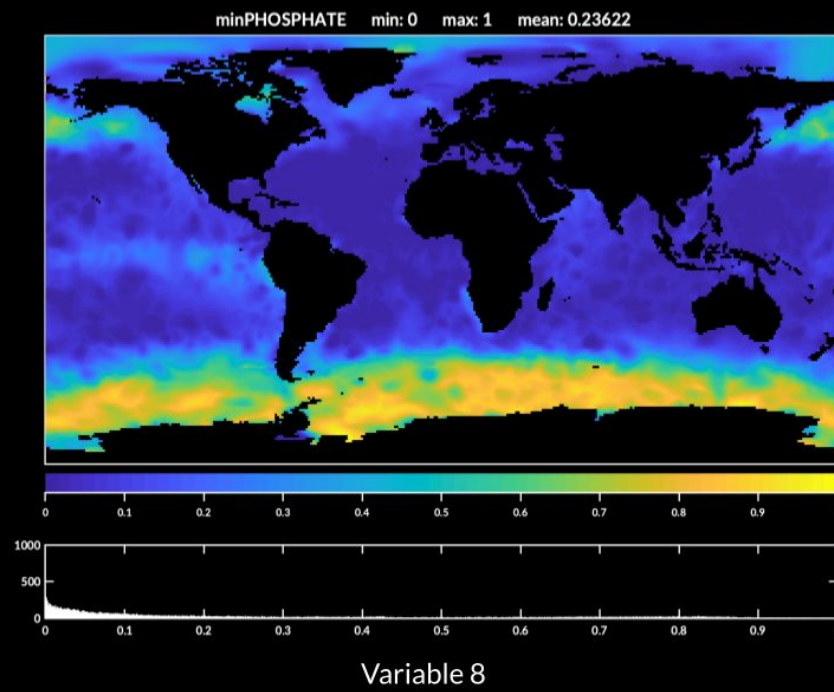
Supplemental Figure 17. Variable 5.



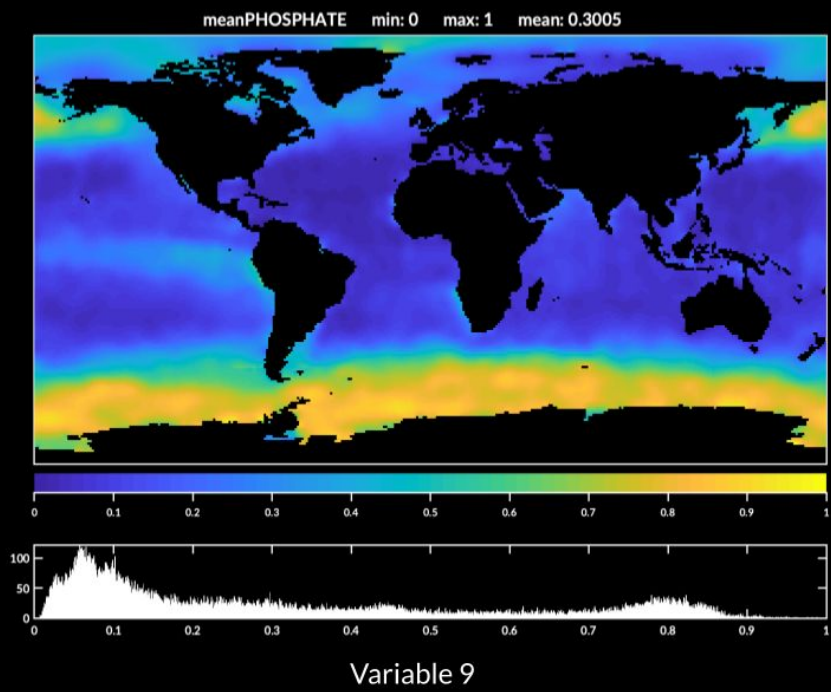
Supplemental Figure 18. Variable 6.



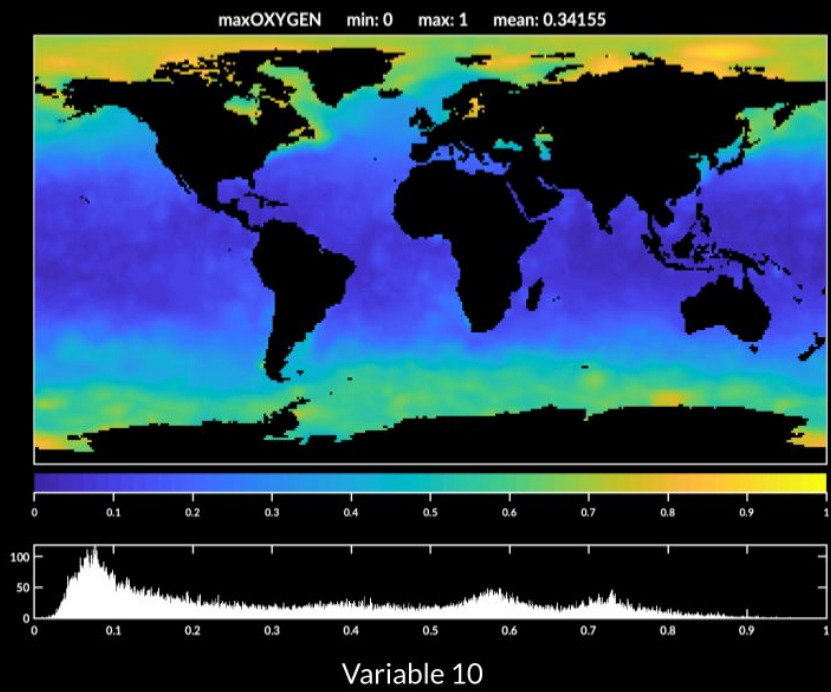
Supplemental Figure 19. Variable 7.



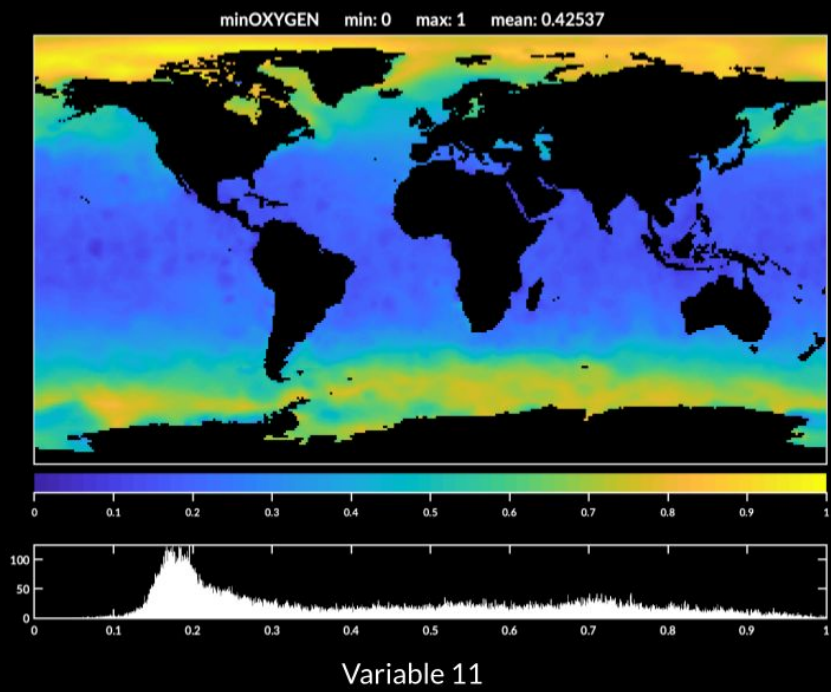
Supplemental Figure 20. Variable 8.



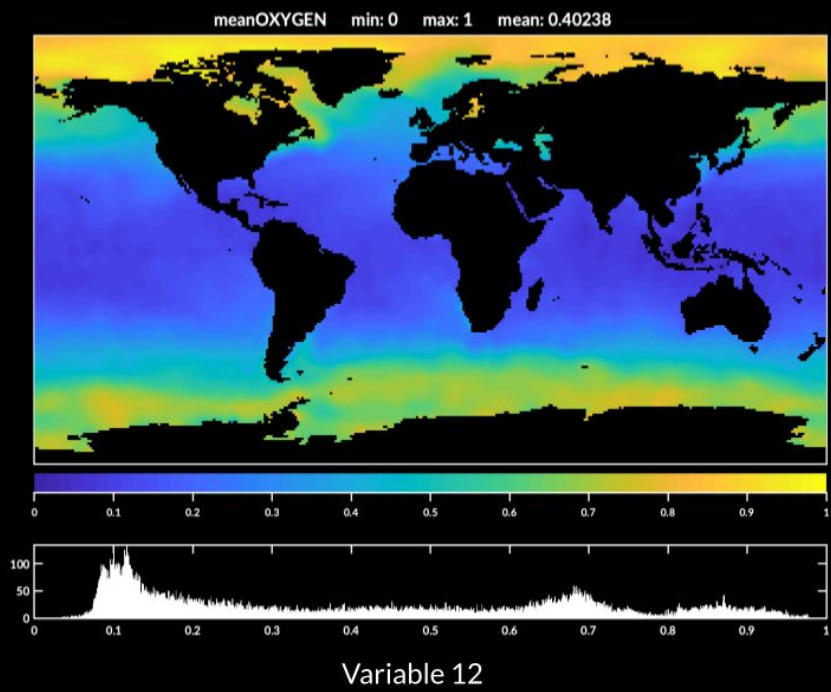
Supplemental Figure 21. Variable 9.



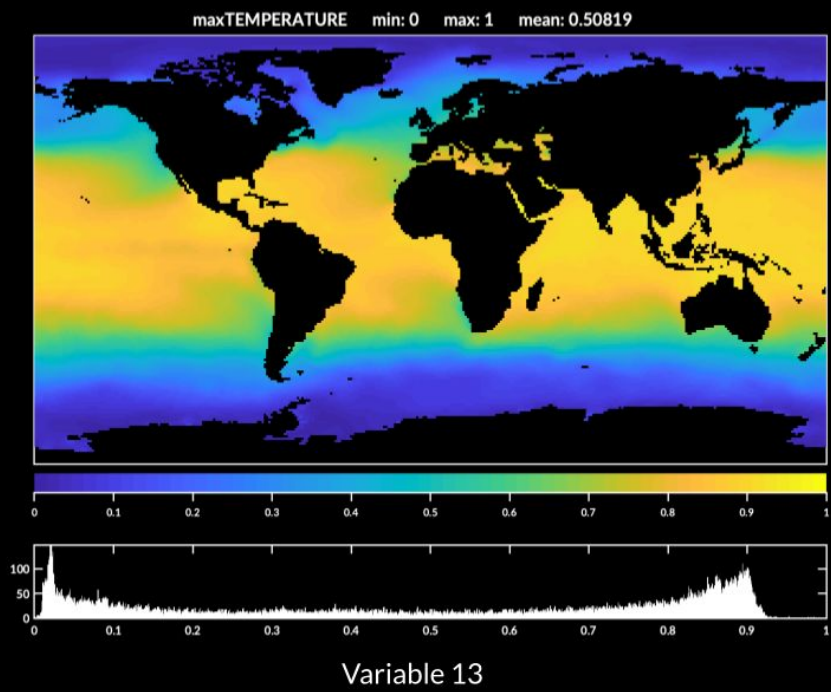
Supplemental Figure 22. Variable 10.



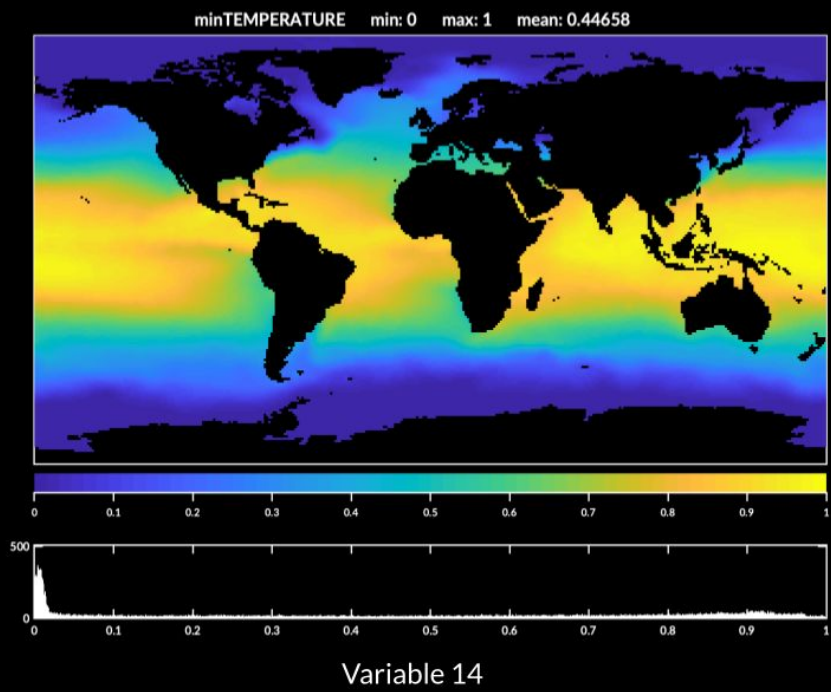
Supplemental Figure 23. Variable 11.



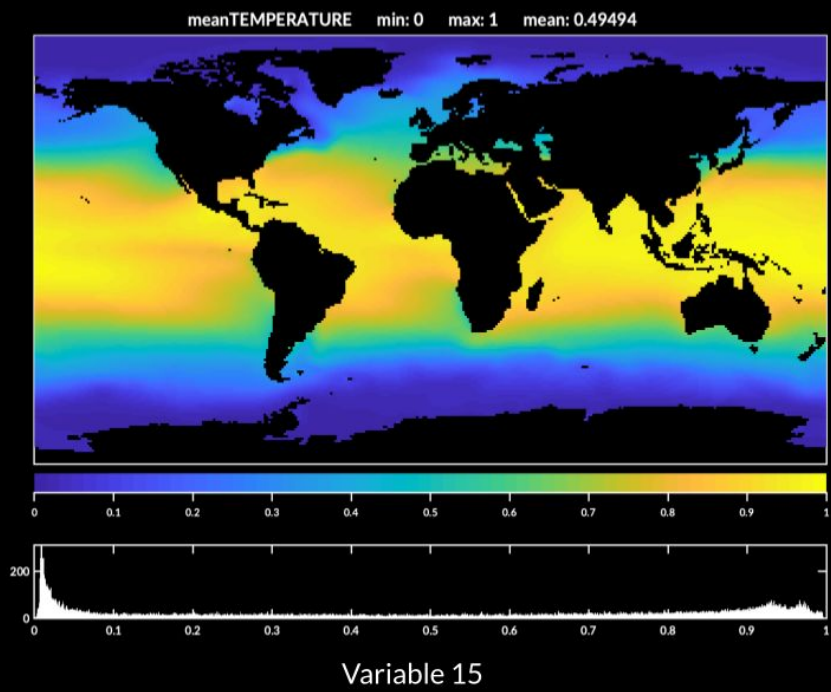
Supplemental Figure 24. Variable 12.



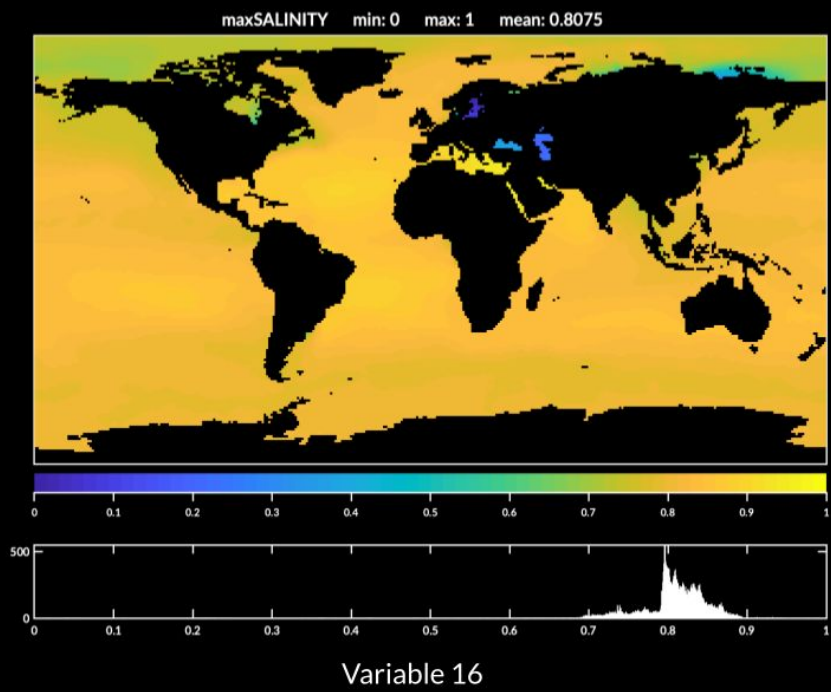
Supplemental Figure 25. Variable 13.



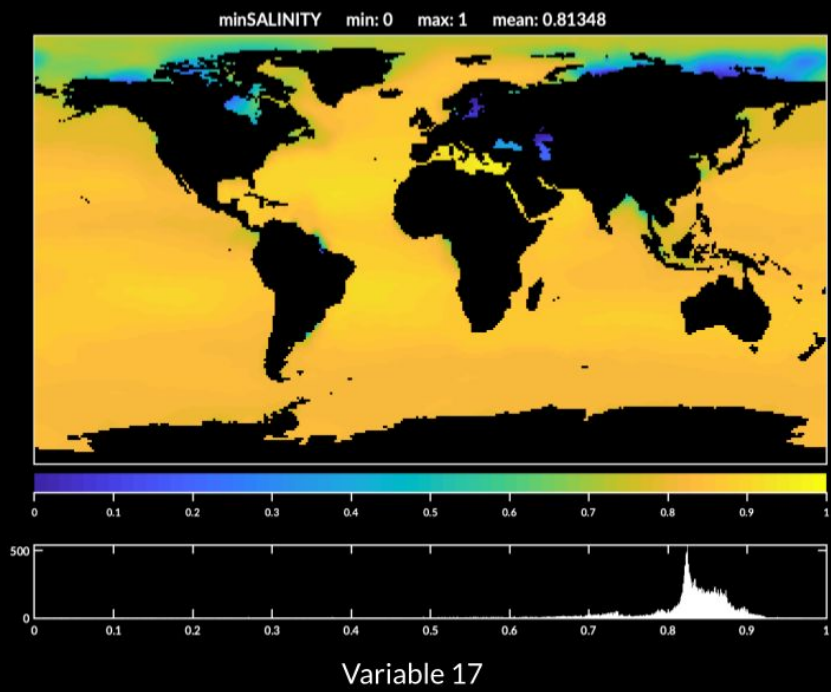
Supplemental Figure 26. Variable 14.



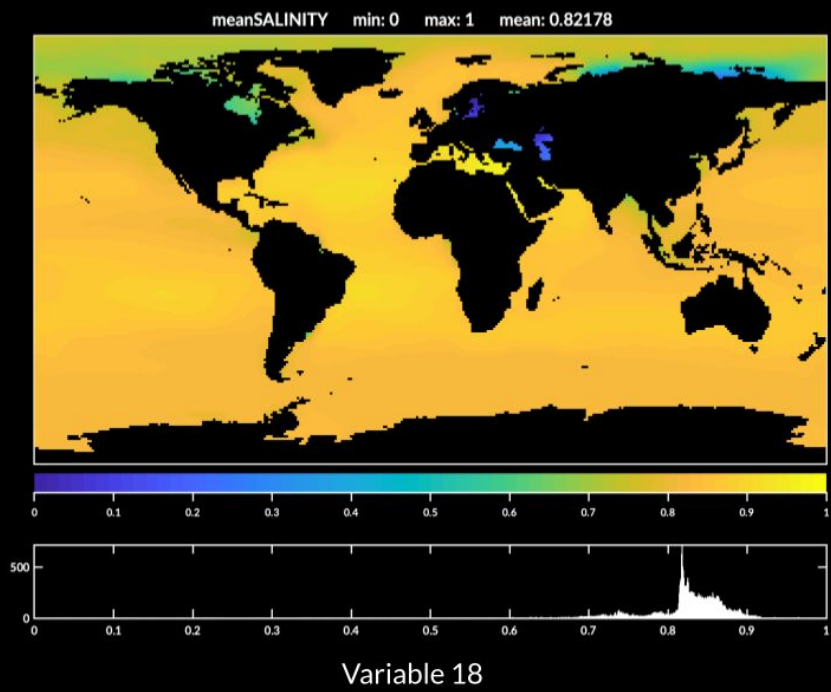
Supplemental Figure 27. Variable 15.



Supplemental Figure 28. Variable 16.



Supplemental Figure 29. Variable 17.



Supplemental Figure 30. Variable 18.