

1 **El Niño detection via unsupervised clustering of Argo**
2 **temperature profiles**

3 **Isabel A. Houghton¹, James D. Wilson^{1,2}**

4 ¹The Data Institute, University of San Francisco, San Francisco, CA 94105

5 ²Department of Mathematics and Statistics, University of San Francisco, San Francisco, CA 94105

6 **Key Points:**

- 7 • Unsupervised clustering based solely on temperature profiles effectively partitions
8 water masses in the Pacific Ocean.
- 9 • The temporal evolution of the clusters reveals spatial oscillations associated with
10 El Niño events.
- 11 • Unsupervised machine learning serves as a flexible and robust approach to anomaly
12 detection in oceanographic data.

Corresponding author: Isabel A. Houghton, izhoughton@gmail.com

Abstract

Variability in the El Niño-Southern Oscillation (ENSO) has global impacts on seasonal temperatures and rainfall. Current detection methods for extreme phases, which occur with irregular periodicity, rely upon sea surface temperature anomalies within a strictly defined geographic region of the Pacific Ocean. However, under changing climate conditions and ocean warming, these historically motivated indicators may not be reliable into the future. In this work, we demonstrate the power of data clustering as a robust, automatic way to detect anomalies in climate patterns. Ocean temperature profiles from Argo floats are partitioned into similar groups utilizing unsupervised machine learning methods. The automatically identified groups of measurements represent spatially coherent, large-scale water masses in the Pacific, despite no inclusion of geospatial information in the clustering task. Further, spatiotemporal dynamics of the clusters are strongly indicative of El Niño events, the east Pacific warming phase of ENSO. The fitting of a cluster model on a collection of ocean profiles identifies changes in the vertical structure of the temperature profiles through reassignment to a different group, concisely capturing physical changes to the water column during an El Niño event, such as thermocline tilting. Clustering proves to be an effective tool for analysis of the irregularly sampled (in space and time) data from Argo floats and may serve as a novel approach for detecting anomalies given the freedom from thresholding decisions. Unsupervised machine learning could be particularly valuable due to its ability to identify patterns in datasets without user-imposed expectations, facilitating further discovery of anomaly indicators.

Plain Language Summary

The climate phenomenon known as El Niño leads to variable temperatures and rainfall amounts around the world and occurs at unpredictable intervals. The most commonly used measurement to determine an El Niño is occurring relies on the difference between the three-month average temperature and the thirty-year average at the surface of the ocean in a rectangular region near the equator. However, as climate changes, these historically defined ways of measuring an El Niño may no longer be helpful. In order to develop a more flexible way to observe an El Niño, we use tools from the field of machine learning. Specifically, temperature measurements in the Pacific Ocean from the surface down to a depth of 1,000 m are grouped automatically (i.e. without pre-defined rules) using machine learning methods. Without using information about the location of the

45 measurements, this process groups measurements that are also close together in space.
46 Changes over time of group assignments closely matches an El Niño happening, and also
47 point to physical changes to that region in the ocean. The automatic grouping of ocean
48 profiles works very well to signal an El Niño and could potentially be a useful tool for
49 future study of data from the ocean.

50 **1 Introduction**

51 The oceans are critical in governing global climate through heat transport and ab-
52 sorption of carbon from the atmosphere (Marshall & Plumb, 2008). Extensive effort is
53 put toward monitoring and predicting the state of the ocean, providing valuable data
54 for daily weather prediction as well as long term understanding of climate variability. The
55 Pacific Ocean, the world’s largest ocean basin, has recurring patterns of variability, most
56 notably as part of the El Niño-Southern Oscillation (ENSO). Due to complex coupling
57 between the ocean and atmosphere, sea surface temperatures and atmospheric winds in
58 the Pacific region interact in a positive feedback loop to produce major oscillations in
59 climate with repercussions at a global scale. An El Niño event, characterized by anoma-
60 lous warming of eastern equatorial Pacific waters, occurs approximately every 3-8 years
61 and, due to global teleconnections, results in varying temperatures and precipitation lev-
62 els around the globe (Wyrтки, 1975; Rasmusson & Carpenter, 1982). The ensuing shift
63 in seasonal temperatures and rainfall leads to droughts and flooding in Africa, Latin Amer-
64 ica, North America, and Southeast Asia. These extreme events have major consequences
65 for human health and economic costs in the billions (Buizer et al., 2000; Iizumi et al.,
66 2014). Despite the importance of forecasting such events, El Niño prediction remains chal-
67 lenging, particularly beyond a six-month horizon, due to the high non-linearity of the
68 system and the relatively unique development of each El Niño event (Timmermann et
69 al., 2018; Dijkstra et al., 2019).

70 The dynamics of the El Niño-Southern Oscillation are associated with a high pres-
71 sure system over the eastern Pacific Ocean and a low pressure system over the western
72 Pacific and Indonesia. This pressure gradient across the Pacific leads to persistent east-
73 erly winds near the equator that drive upwelling along the eastern Pacific coasts, lead-
74 ing to cooler surface temperatures and a tilted thermocline. During an El Niño event,
75 the pressure gradient driven atmospheric circulation decreases, reducing upwelling along

76 the eastern Pacific, enhancing sea surface temperatures and deepening of the thermo-
77 cline in that region (Wang et al., 2000; Meinen & McPhaden, 2000).

78 Current El Niño detection relies on sea surface temperature anomalies within a specif-
79 ically designated region (Niño 3.4, defined from 5°S to 5°N and 170°W to 120°W) in the
80 equatorial Pacific. Extensive study of historical patterns have identified this region as
81 the dominant location of the coupled ocean-atmosphere interactions (Bamston et al., 1997).
82 The exclusive consideration of surface measurements in a small geographic location po-
83 tentially disregards indicators in other regions of the Pacific Ocean basin and in subsur-
84 face variation of the vertical structure. As a consequence of a significant warming trend
85 in the Niño 3.4 region since 1950, the thirty-year period against which temperature val-
86 ues are compared is updated on a five-year basis (NOAA National Centers for Environ-
87 mental Prediction, 2020). In the context of global climate change and ocean warming,
88 these updates will likely continue to be necessary (Ashok et al., 2007; Yeh et al., 2009).
89 Therefore, methods for El Niño detection incorporating large horizontal and vertical scales
90 and utilizing in situ data without empirical thresholds are of particular value (Yang &
91 Wang, 2009).

92 In situ measurements of the ocean are valuable sources for subsurface observations
93 as well as for model validation and improvement, particularly in a changing climate. In
94 situ instruments have begun collecting increasing amounts of data, thus methods for ef-
95 fective analysis are critical for data utilization and could provide new approaches to ocean
96 observation and prediction. To date, the Argo program (Riser et al., 2016) has massively
97 increased the extent of in situ measurements of the ocean with profiling floats. Those
98 direct measurements have provided insight into ENSO dynamics, specifically allowing
99 detection of changes in the distribution of ocean heat content due to tilting of the ther-
100 mocline (Roemmich & Gilson, 2011; Johnson & Birnbaum, 2017) and adjustment of the
101 Equatorial Pacific Thermocline, a layer of low vertical stratification below the pycnocline
102 (Johnson & Birnbaum, 2016). In situ measurements do come with additional challenges
103 over uniform model data, particularly in terms of spatial and temporal sparsity (rela-
104 tive to model grid cells) and nonuniform sampling. As a result, development of novel meth-
105 ods for data utilization may prove particularly useful with the growing deluge of data
106 becoming available.

107 Unsupervised machine learning methods for clustering data provide an effective and
108 robust approach for partitioning complex data, particularly adaptable to the spatial and
109 temporal irregularity of many in situ ocean observations. Additionally, clustering can
110 reveal patterns or similarities in a dataset while avoiding biased expectations of what
111 patterns should exist (e.g., thresholds derived from prior assumptions of the system). Pre-
112 vious work developed a profile classification model using clustering and considered un-
113 supervised clustering of temperature profile measurements in the Atlantic and South-
114 ern Oceans (Maze, Mercier, & Cabanes, 2017; Maze, Mercier, Fablet, et al., 2017; Jones
115 et al., 2019) and found groupings consistent with known oceanic water masses. In this
116 work, we analyze measurements in the Pacific Ocean basin and consider the temporal
117 evolution of the clustered data. The openly-available dataset of ocean temperature pro-
118 files from the Argo program is analyzed with unsupervised machine learning methods
119 to reveal novel El Niño indicators free from user-imposed decisions. We find that tem-
120 poral dynamics in the spatial location of cluster assignments are strongly correlated with
121 the current leading metric for El Niño occurrence. The unsupervised methods success-
122 fully partition the temperature profiles into physically meaningful groups and the vari-
123 ation over time identifies changes in both thermocline depth and sea surface tempera-
124 tures, key physics associated with ENSO. The data and analysis methods are described
125 in the following section. Section 3 describes the patterns identified by the clustering al-
126 gorithm and section 4 discusses their relationship to current oceanographic understand-
127 ing. Finally, section 5 summarizes the utility of unsupervised methods for analyzing ocean-
128 ographic data as illustrated by effective ENSO detection and highlights future directions.

129 **2 Data and Methods**

130 Temperature profiles in the Pacific Ocean acquired by the Argo project (Riser et
131 al., 2016) were reduced to a lower-dimensional embedding using principal component anal-
132 ysis (PCA) and then grouped via k-means clustering, an unsupervised clustering method.
133 The evolution of the spatial patterns of measurements assigned to each cluster were then
134 considered over a thirteen-year time period (2006–2019). Measurements within each three-
135 month period during this time period were aggregated for analysis. Oscillations in the
136 spatial extent of clusters were compared to an indicator of El Niño. A description of the
137 Argo temperature dataset, dimensionality reduction and clustering methods, and meth-

138 ods used to evaluate the spatial evolution of our clusters relative to an existing El Niño-
139 Southern Oscillation indicator over time are included below.

140 **2.1 Argo Float Dataset**

141 The Argo program was initiated in the early 2000's and consists of a global array
142 of free-drifting profiling floats that have served to massively expand our global ocean ob-
143 serving network. Each profiling float in the array measures the vertical structure of tem-
144 perature and salinity in the ocean, with newer profiling floats also measuring bio-optical
145 traits and biogeochemical properties. Currently, nearly 4,000 individual profiling floats
146 are deployed, each acquiring vertical profile measurements to a depth of approximately
147 2,000 m every ten days. Collected data is then made publicly available in near real-time.
148 The free-floating nature of the instruments leads to a global array of sensors distributed
149 at roughly every three degrees (~ 300 km), with dynamically changing positions over time.
150 Argo is the leading source of global subsurface data, particularly for use in ocean data
151 assimilation and model reanalysis (Riser et al., 2016).

152 Argo profiling float measurements of temperature were acquired in October 2019
153 in the Pacific Ocean basin between 30°S and 50°N from January 2006 to September 2019
154 (Argo, 2019). The start date was chosen to achieve a sufficiently long window to observe
155 several El Niño events while utilizing a similar number of measurements at a given time
156 (see supplementary Figure S1). Prior to 2006, the number of profile measurements avail-
157 able becomes sparser. The latitude range was chosen to initially focus on the low- to mid-
158 latitudes of the Pacific basin. Each measurement had an associated latitude, longitude,
159 and acquisition timestamp. Only delayed-mode data were used. All temperature pro-
160 files containing missing data, insufficient data points, or nonphysical values were removed
161 as defined below. This corresponded to profiles with fewer than 50 data points, the ini-
162 tial data point more than 25 dbar from the surface, the final data point less than 1,000
163 dbar, or temperature values less than -5°C . Temperature values in the remaining pro-
164 files were linearly interpolated onto a uniform grid with 5 dbar spacing from 5 dbar down
165 to 1,000 dbar, the same as Jones et al. (Jones et al., 2019). Data was only stored down
166 to 1,000 dbar despite measurements extending down to approximately 2,000 dbar due
167 to the majority of temperature variability of interest associated with the thermocline and
168 occurring in the upper 1,000 dbar (Yang & Wang, 2009). This yielded a set of approx-
169 imately 560,000 temperature profiles consisting of 200 data points each for the thirteen

170 year time span that were subsequently assigned to clusters. The count of profiles increased
171 from approximately 35,000 per year in 2006 to approximately 70,000 in 2019 (see sup-
172 plementary Figure S1).

173 2.2 Dimensionality Reduction and Clustering

174 A critical first step toward effective clustering for a high-dimensional variable is di-
175 mensionality reduction (Beyer et al., 1998). Effective dimensionality reduction casts a
176 given sample with many features into a lower-dimensional space where a distance met-
177 ric between two samples reasonably captures differences within the dataset. For the tem-
178 perature profiles consisting of hundreds of data points over a uniform depth grid, cal-
179 culating a point-wise difference between each profile would not fully capture critical dif-
180 ferences between profiles, such as the shape of the temperature profile with depth (e.g.
181 thermocline location).

182 In this work, principal component analysis (PCA) was applied utilizing the *scikit-*
183 *learn* machine learning library for Python (Pedregosa et al., 2011). This algorithm im-
184 plements linear dimensionality reduction using singular value decomposition of the data
185 to project each sample into a lower dimensional space of linearly uncorrelated (orthog-
186 onal) basis functions, termed principal components (Shlens, 2003). The first principal
187 component accounts for the largest possible variance in the data, and each subsequent
188 component attempts to further maximally account for variance under the constraint of
189 orthogonality to preceding components. Thus, one can specify the desired variance to
190 account for in the data and the number of components to describe that variance between
191 samples will be retained. PCA was applied to the 200-data-point profiles to capture 99.9%
192 of the variance with 17 principal components. Across all of the profiles, each depth level
193 was mean-centered but not scaled, therefore the lower depth levels with less tempera-
194 ture variability contributed less to the final representation. Seventeen components was
195 notably higher than previous work by Jones et al. (2019), which only required six com-
196 ponents for 99.9% of the variance, likely due to the strong vertical coherence of the South-
197 ern Ocean (Karsten & Marshall, 2002) as well as Maze et al. (2017b), which prescribed
198 retaining 11 principal components for 99.88% of the variance. Beyond the first few, the
199 principal components had very small magnitudes which is perhaps why an additional six
200 principal components were necessary to account for 99.9% of the data variance here com-
201 pared to the 99.88% used in Maze et al. (2017b).

202 With dimensionality reduction applied, properties such as Euclidean distance be-
 203 tween each new feature (i.e., principal component) become notably more effective at de-
 204 scribing sample differences (Beyer et al., 1998). Clustering methods were next applied
 205 with the goal of grouping the profiles solely based on differences in temperature struc-
 206 ture without any geospatial information or external constraints applied. A wide variety
 207 of clustering methods exist with different advantages and levels of complexity (Xu & Tian,
 208 2015). While exploration of the different clustering outcomes from the variety of meth-
 209 ods (e.g. spectral clustering, hierarchical models) would potentially reveal interesting in-
 210 sights, the primary goal of this study was to find a straightforward approach to assign
 211 temperature profiles to groups. Previous work utilized Gaussian mixture modeling (GMM),
 212 which aims to fit the data as a linear combination of multidimensional Gaussian distri-
 213 butions. In this work, k-means clustering, a widely utilized and efficient approach in a
 214 variety of applications (Jain, 2010), was chosen, primarily due to its computational ef-
 215 ficiency and straightforward implementation. Results from k-means were compared with
 216 GMM (see supplement).

217 Given a set of samples $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$, where each sample is represented by a d -
 218 dimensional vector, the k-means clustering algorithm aims to partition the n samples
 219 into k clusters, $\mathbf{C}=\{C_1, C_2, \dots, C_k\}$, with the objective of minimizing the within-cluster
 220 sum of squares (WCSS). In particular, let μ_ℓ be the mean of the data within the ℓ th clus-
 221 ter, C_ℓ . The k-means algorithm seeks to identify the partition, \mathbf{C} , that minimizes

$$222 \quad WCSS = \sum_{\ell=1}^k \sum_{\mathbf{x} \in C_\ell} \|\mathbf{x} - \mu_\ell\|^2. \quad (1)$$

223 The embeddings of the temperature profiles produced by PCA were clustered fol-
 224 lowing the *scikit-learn* implementation of the k-means clustering task to assign each pro-
 225 file measurement to a cluster.

226 One limitation of k-means clustering lies in the required choice of number of clus-
 227 ters, k , to create. However, due to the efficiency of implementation of the algorithm, a
 228 range of cluster counts can be tested and cluster characteristics can be analyzed to as-
 229 sess optimal cluster count. A common strategy to assess the cohesion of clusters in a par-
 230 tition, i.e. how similar every object is to its cluster, is to measure the average silhouette
 231 score of the cluster assignment (Rousseeuw, 1987).

To obtain a silhouette score, for each data point $i \in C_\ell$, the mean distance between i and all other data points in the same cluster is given by:

$$a(i) = \frac{1}{|C_\ell| - 1} \sum_{j \in C_\ell, i \neq j} d(i, j) \quad (2)$$

where $d(i, j)$ is the distance between cluster points i and j in the cluster C_ℓ , and $|C_\ell|$ denotes the number of data points in cluster ℓ . The dissimilarity of point $i \in C_\ell$ to other clusters is then defined by:

$$b(i) = \min_{k \neq \ell} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j) \quad (3)$$

where the cluster to which sample i is closest, but not assigned, is used (indicated by the *min* operator). Combining the similarity of a sample to its assigned cluster ($a(i)$) and dissimilarity to the nearest cluster to which it is not assigned ($b(i)$), yields a silhouette score, s , defined as:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad (4)$$

which can then be aggregated for all partitioned points. To assess the cohesion of a partition, \mathbf{C} , we measure the average silhouette score across all data points. An optimal silhouette score of 1.0 indicates the sample is a large distance to non-assigned clusters and small distance to other samples in the assigned cluster. A score of -1.0 indicates the sample is closer to another cluster than its own and a score of 0 indicates the object is on the border of two natural clusters. The global silhouette score can be calculated for varying cluster counts, ideally encountering a cluster count, k , that maximizes the global silhouette score. The silhouette score was taken into account with physical intuition regarding the Pacific Ocean in order to find an optimal cluster count that maximizes uniqueness of data in the clusters with sufficient clusters to describe variability in the Pacific. Specifically, inspection of the unique water masses in the Pacific Ocean (Emery, 2008) indicated likely more than three clusters (the value of k with the global maximum silhouette score, see below and Figure 1) would be useful to capture the variability given the presence of several upper water masses in the Pacific (e.g. Equatorial, Eastern and Western Central waters, Northern and Southern Eastern Transition waters) overlapping with intermediate waters (e.g. North Pacific Intermediate Water, California Intermediate Water, and Antarctic Intermediate Water).

K-means clustering was found to be effective at partitioning, reproducible, and highly computationally efficient. The silhouette score for cluster counts ranging from 3 to 10

263 exhibited a stable point at $k = 7$ (Figure 1), indicating partitioning at that granularity
264 aligned with separations in the data. Seven clusters were chosen in order to balance ob-
265 taining a reasonable number of clusters with improvement seen in the silhouette score.
266 While choice of k did involve decision making in an otherwise unsupervised process, vari-
267 ation of cluster count did not fundamentally alter the partitioning, but rather led to a
268 coarsening (for fewer clusters) or refining (for more clusters) of the divisions along sim-
269 ilar lines (see supplementary Figure S3). Following selection of an appropriate k , data
270 across all time (2006-2019) were simultaneously clustered and the assigned cluster iden-
271 tity was used for subsequent analysis. Alternatively, temperature profiles could be di-
272 vided into shorter time periods and then subsequently clustered (not shown). However,
273 simultaneous clustering across all time yielded similar partitions and provided a more
274 consistent approach, particularly given the free-floating, intermittent nature of the mea-
275 surements in contrast to a fixed set of sampling locations.

276 Because the k-means clustering algorithm is randomly initialized and can converge
277 on a local (rather than global) minimum, repeatability of the clustering assignment was
278 quantified with an adjusted Rand index measuring the similarity between two different
279 groupings, adjusted for random chance of assignment (Rand, 1971). An index of 1.0 in-
280 dicates exactly identical clustering, regardless of specific label changes (i.e. a cluster la-
281 belled #1 in one partitioning can be labelled cluster #4 in a subsequent partitioning but
282 have the same members). The adjusted Rand index was calculated before analysis was
283 carried out to confirm that repeated clustering would yield similar results. Repetition
284 of the clustering produced very similar results such that the same profiles were consis-
285 tently grouped together. Ten repeated clusterings produced an average adjusted Rand
286 index of 0.997, indicating high repeatability of the analysis.

287 **2.3 El Niño-Southern Oscillation Indicator**

288 The current leading diagnostic metric of El Niño-Southern Oscillation state uti-
289 lized by the National Oceanic and Atmospheric Administration (NOAA) relies on the
290 sea surface temperature anomaly within the rectangular Niño 3.4 region of the Pacific
291 defined from 5°S to 5°N and 170°W to 120°W (NOAA National Centers for Environ-
292 mental Prediction, 2020). The three-month running mean of the anomaly from the most
293 recent 30-year historical period (updated every five years) in this region is termed the
294 Oceanic Niño Index (ONI). This index must exceed $\pm 0.5^\circ\text{C}$ for at least five consecutive

295 overlapping three-month periods to classify the period as a full-fledged El Niño (+0.5°C)
 296 or La Niña (-0.5°C) (NOAA National Centers for Environmental Prediction, 2020). ONI
 297 values were obtained from NOAA (NOAA National Centers for Environmental Predic-
 298 tion, 2020) and used directly for comparison.

299 **2.4 Spatio-temporal Cluster Analysis**

300 Following the clustering of temperature measurements without any associated tem-
 301 poral or geospatial information, the locations of measurements assigned to each cluster
 302 were analyzed over time and compared to historic El Niño events, utilizing the ONI as
 303 a ground truth on the historic presence or absence of an event. All profile measurements
 304 occurring in a 90 day window were aggregated into a single timestep with the window
 305 shifting by 30 days for each subsequent timestep, providing statistics representing a three-
 306 month running mean for comparison with the ONI values. For each cluster (ℓ) and for
 307 each timestep (t), the anomaly from the average longitude of the cluster over the full 13-
 308 year period ($\lambda_{\ell,2006-2019}$) was then considered. However, there could be many measure-
 309 ments in a given area (because the profiling floats are not enforced to be uniformly dis-
 310 tributed) skewing the mean toward that region. To effectively capture how far a given
 311 cluster extends east or west beyond its $\lambda_{\ell,2006-2019}$ position at a given timestep, all unique
 312 longitudes of measurements within the cluster were aggregated, essentially indicating the
 313 spatial coverage of that cluster. In practice, all measurements in the same 0.5 degree lon-
 314 gitude bin were considered as a single unique longitude. The average of this unique set
 315 of longitudes was calculated ($\lambda_{\ell,t}$) and then differenced from the $\lambda_{\ell,2006-2019}$ position
 316 for this cluster to determine the anomaly in the average of the unique longitudes for that
 317 cluster at that timestep. This method minimized the importance of several measurements
 318 at the same longitude (but potentially different latitude) and highlighted oscillations in
 319 the zonal extent of a cluster.

320 **3 Results**

321 **3.1 Clustering**

322 Each group produced by the clustering algorithm contained profiles with relatively
 323 similar vertical structure and surface temperature values indicated by the uniqueness of
 324 the average temperature profile of each cluster and the standard deviation within the

325 group relative to variation between groups (Figure 2). The unsupervised clustering method
 326 was able to detect differences and partition profiles with similar surface temperatures
 327 but unique vertical structures (e.g. clusters 0 and 5), as well as similar vertical struc-
 328 tures but shifted temperatures (e.g. clusters 2 and 5), a complex task to achieve with
 329 hard-coded selection rules. Each measurement assigned to a cluster also had an associ-
 330 ated latitude and longitude allowing visualization of clusters in geographic space. A map
 331 of all measurement locations, each colored by its corresponding cluster assignment, il-
 332 lustrates the spatial coherency of each cluster, with few outliers and minimal spatial over-
 333 lap between clusters (Figure 3a). This spatial coherency was similar to previous anal-
 334 yses by Maze, Mercier, Fablet, et al. (2017) and Jones et al. (2019), despite utilization
 335 of a different clustering method (k-means versus Gaussian mixture model). Notably, when
 336 only sea surface temperature (i.e. the uppermost measurement by the profiling float) was
 337 used for clustering (Figure 3b), the clusters were significantly less spatially well-defined
 338 with a scattered overlap of measurements belonging to different groups, indicating the
 339 subsurface structure of the temperature profile was critical in partitioning.

340 **3.2 Spatio-Temporal Dynamics**

341 Assignment of measurements to clusters from three-month time periods exhibited
 342 clear spatial oscillations correlated with the Oceanic Niño Index. Oscillations were pri-
 343 marily observed in clusters with measurements at lower latitudes (see Figure 4 and sup-
 344 plementary video). Figure 4 revealed a noticeable change in clustering assignments which
 345 closely matched El Niño events.

346 **3.2.1 Niño 3.4 Region**

347 For direct comparison with the current region considered for diagnosis of El Niño
 348 conditions, clustering of measurements in the constrained geographic region of Niño 3.4
 349 (N3.4) was considered first. The cluster assignments, rather than the traditional surface
 350 temperature values, were analyzed. Two groups primarily populated the N3.4 region over
 351 the thirteen years, a low latitude western group (cluster 5, teal) and a low latitude east-
 352 ern group (cluster 2, orange). The two groups occupied unique spatial regions with an
 353 east-west division. Qualitatively, the division oscillated east and west irregularly, in syn-
 354 chrony with the ONI (inner boxed regions, Figure 4). During neutral ENSO periods, the
 355 N3.4 region was approximately evenly divided between one group in the western half and

356 one group in the eastern half (Figure 4a,c). During a positive ONI anomaly (El Niño event),
 357 the western cluster distinctly shifted eastward to occupy the majority of the N3.4 region
 358 (Figure 4b,d). Following an event, as the ONI rapidly returned to neutral levels, the west-
 359 ern cluster shifted back to its original balance partially occupying the N3.4 region along
 360 with eastern cluster measurements. The shifting of the spatial locations of measurements
 361 assigned to a group is quantified by the anomaly in the average unique longitudes of mea-
 362 surements in the eastern cluster (Figure 5a). By the anomaly of the $\lambda_{2,t}$ positions for
 363 the eastern cluster at each time step from that cluster's $\lambda_{2,2006-2019}$ position as defined
 364 above, the average unique longitudinal position of measurements in cluster 2 was con-
 365 sistenty farther east (positive longitudinal anomaly) during periods above the El Niño
 366 threshold, and near average or farther west during other periods.

367 **3.2.2 Tropical Pacific Region**

368 Temporal dynamics of cluster assignments in the entire tropical Pacific region span-
 369 ning $\pm 23.4^\circ$ latitude indicated additional larger-scale patterns. The tropics were pri-
 370 marily populated by three groups: one group (cluster 2, orange) in the eastern Pacific
 371 spanning the tropical latitudes, a second group in the western Pacific confined to lower
 372 latitudes (cluster 5, teal), and a third group (cluster 0, maroon) also in the western Pa-
 373 cific to the north and south of the second group (Figure 3a). During an elevated ONI
 374 period, the eastern cluster that had shifted farther east at very low latitudes (N3.4 re-
 375 gion), simultaneously significantly expanded its extent westward at slightly northern lat-
 376 itudes, leading to the presence of measurements assigned to this cluster all the way across
 377 the Pacific in a narrow band around 10°N (Figure 4b,d). This phenomenon exhibited
 378 itself during every El Niño event during the time period assessed (2006-2019). This os-
 379 cillation was quantified with the anomalous average unique longitude ($\lambda_{2,t} - \lambda_{2,2006-2019}$)
 380 of the eastern cluster (Figure 5b), and had a Pearson correlation coefficient with the ONI
 381 of -0.75 and a peak cross-correlation with zero time lag.

382 **4 Discussion**

383 The ocean is composed of a distribution of water masses with unique temperature
 384 and salinity characteristics that can be related to the region of water mass formation (Emery,
 385 2008). These water masses typically have both a horizontal and vertical extent. There-
 386 fore, a profile measurement down to 1,000 dbar would likely sample multiple water masses,

387 indicated by temperature and salinity variability over depth in the profile. This layer-
388 ing of unique water masses with variable horizontal extents results in the high variabil-
389 ity seen in temperature profiles. However, temperature profiles obtained physically prox-
390 imate are likely sampling the same set of water masses and therefore likely to exhibit sim-
391 ilar structure. The effective clustering of similarly structured temperature profiles in turn
392 led measurements within a given cluster to be spatially proximate, as seen in Figure 3.
393 Further, the clusters identified align well with traditional water masses and the pattern
394 of overlap of different depth water masses. Specifically, cluster 3 (yellow) aligns well with
395 the boundaries of Pacific Subarctic Water (PSUW) (Figure 3). Cluster 4 (light green)
396 appears to represent the regions occupied by the North and South Pacific Subtropical
397 Mode Waters (NPSMW, SPSMW). The path of the Kuroshio Current aligns well with
398 the boundary between clusters 4 (light green) and 1 (dark orange) while the path of the
399 Oyashio Current appears to align well with the boundary between clusters 3 (yellow) and
400 1 (dark orange). The boundaries on the meridional extent of the Pacific Equatorial Wa-
401 ter relative to the Western North Pacific Central Water (WNPCW) and the Western South
402 Pacific Central Water (WSPCW) in the western equatorial Pacific can also be seen as
403 the boundary between the clusters 5 (teal) and 0 (red) (Figure 3). Intermediate depth
404 waters are also formed off the coast of California in the northern hemisphere and off the
405 coast of South America in the southern hemisphere, termed California Intermediate Wa-
406 ter (CIW) and East Southern Pacific Intermediate Water (ESPIW) (Emery, 2008), re-
407 spectively, as a consequence of coastal upwelling. Both of these water masses appear clus-
408 tered together in cluster 1 (dark orange) near the sites of coastal upwelling. Meanwhile,
409 cluster 2 (light orange) appears to represent water masses related to equatorial dynam-
410 ics and Eastern South Pacific Central Water/Eastern North Pacific Central Water. The
411 Pacific Equatorial Water (PEW) forms a notable band of water at low latitudes (Emery,
412 2008). This region was also partitioned by the clustering task, and was divided into an
413 eastern and western cluster at low latitudes. This east-west division of the PEW was due
414 to the variable thermocline depth and surface temperature across the Pacific, with suf-
415 ficiently high variability relative to the rest of the profiles for the algorithm to identify
416 two unique clusters (Figure 2). Interestingly, even very few partitions (e.g. $k = 3$) still
417 divided the PEW, indicating the east-west variability was relatively dominant (see sup-
418 plement, Figure S3). Despite thermocline depth and surface temperature varying con-
419 tinuously across the Pacific, the partitioning divided the PEW at a consistent thresh-

old, with the location of that division found to be particularly relevant in terms of temporal variability.

The switching of cluster assignment in the regions of interest signaled a physical change to the water column indicated by the differences in temperature profiles in the two dominant oscillating clusters (i.e. orange and teal profiles in Figure 2). The eastern cluster is characterized by cooler surface temperatures and a shallower thermocline, therefore a shift of that cluster out of the N3.4 region aligns with the positive ONI temperature anomaly. All anomalies in the average of the unique longitudes within the eastern cluster ($\lambda_{2,t} - \lambda_{2,2006-2019}$) beyond one standard deviation occur simultaneously with an El Niño event, and only the major event in 2015-2016 exceeds two standard deviations (Figure 5). At the surface, the profiles in the western cluster (5) have warmer temperatures than profiles in the eastern cluster (2). In terms of vertical structure, the thermocline is deeper in the western cluster and shallower in the eastern cluster. Thus, during neutral conditions, the east-west division in the two clusters corresponds to a tilted thermocline and colder surface temperatures in the east. During an El Niño, the western cluster extends farther eastward at the equator, indicating warmer surface temperatures and a deeper thermocline than under neutral conditions, consistent with physical understanding of ENSO dynamics (Meinen & McPhaden, 2000). Additionally, the eastern cluster extends far westward in a band north of the western cluster, leading to a north-south gradient in cluster identity and accompanying north-south surface temperature gradient and thermocline tilt that is unique to periods with an elevated Oceanic Niño Index. The spatial extent of the clusters thus provided a concise method for observation of oscillations characteristic of Kelvin and Rossby wave-driven ENSO dynamics (Battisti, 1989; Kim & Kim, 2002). The ability to compare the general characteristics of profiles in each group produced by the clustering provided a concise way to identify complex shifts in water column structure over time and clearly identify anomalous periods.

Unsupervised clustering provided a robust way to delineate regions with distinct water masses without imposing thresholds or arbitrary latitude or longitude limits. Additionally, the spatial locations of measurements within a cluster evolved over time, and relating back to the original temperature profiles in a given cluster indicated the physical dynamics at work, such as a shift in thermocline depth.

5 Conclusions

452
453 Approximately 560,000 temperature profiles in the Pacific Ocean taken from 2006-
454 2019 were partitioned into seven groups via the k-means clustering method. Analysis of
455 all measurement assignments illustrate spatially coherent patterns associated with known
456 water masses of the Pacific despite no inclusion of geospatial information in the cluster-
457 ing decision. Cluster assignments over time oscillate in spatial extent, particularly at lower
458 latitudes. These oscillations are strongly correlated with the Oceanic Niño Index, the
459 broadly utilized indicator of an El Niño event. The representative profiles of each cluster
460 correspond to the current understanding of oceanic dynamics, particularly the shift
461 in sea surface temperature and thermocline depth as a result of reduced eastern Pacific
462 upwelling during El Niño events. Despite the difficult task of uniformly sampling a mas-
463 sive extent of the worlds oceans with free-drifting devices, Argo sensors are gathering
464 sufficient data to observe oscillations in oceanic dynamics over relatively short time pe-
465 riods (i.e. three months) at relatively high resolution (3-5 degrees), indicating the un-
466 paralleled value of the ever-increasing observing network and the real-time data distri-
467 bution.

468 While clustering methods have been applied across a variety of fields, utilization
469 within ocean and climate sciences remains limited (Karpatne et al., 2019). However, as
470 climate change continues and potentially accelerates (IPCC, 2019), identifying robust
471 methods to identify patterns and anomalies within climate and environmental data could
472 prove invaluable as historic means continuously shift. In the context of climate models
473 in a changing climate, this objective approach could further serve to account for biases
474 in ENSO representation. Unsupervised methods such as clustering and other complex
475 network theory approaches (e.g. anomaly detection on a graph) provide an automated
476 approach to segmentation and analysis driven by statistics of the dataset rather than po-
477 tentially imposing biases toward expected, but not necessarily fully representative, pat-
478 terns.

479 Altogether, unsupervised machine learning techniques prove to be a highly effec-
480 tive approach for analyzing Argo data and gaining physical insights into the system.

481 **Acknowledgments**

482 These data were collected and made freely available by the International Argo Program
 483 and the national programs that contribute to it; <http://www.argo.ucsd.edu>, <http://argo.jcommops.org>.
 484 The Argo Program is part of the Global Ocean Observing System. The ONI was pro-
 485 vided by the National Oceanic and Atmospheric Administration National Centers for
 486 Environmental Prediction; [https://origin.cpc.ncep.noaa.gov/products/analysis_ moni-
 487 toring/ensostuff/ONI_v5.php](https://origin.cpc.ncep.noaa.gov/products/analysis_monitoring/ensostuff/ONI_v5.php).

488 JDW was partially funded by the National Science Foundation grant NSF DMS-
 489 1830547. IAH was supported by The Data Institute at University of San Francisco.

490 **References**

- 491 Argo. (2019). Argo float data and metadata from Global Data Assembly Center
 492 (Argo GDAC) - Snapshot of Argo GDAC of October 8st 2019. *SEANOE*. Re-
 493 trieved from <https://www.seanoe.org/data/00311/42182/#67548> doi: 10
 494 .17882/42182
- 495 Ashok, K., Behera, S. K., Rao, S. A., Weng, H., & Yamagata, T. (2007). El Niño
 496 Modoki and its possible teleconnection. *Journal of Geophysical Research:*
 497 *Oceans*. doi: 10.1029/2006JC003798
- 498 Bamston, A. G., Chelliah, M., & Goldenberg, S. B. (1997). Documentation of a
 499 highly enso-related sst region in the equatorial pacific: Research note. *Atmo-*
 500 *sphere - Ocean*. doi: 10.1080/07055900.1997.9649597
- 501 Battisti, D. S. (1989). On the Role of Off-Equatorial Oceanic Rossby Waves dur-
 502 ing ENSO. *Journal of Physical Oceanography*. doi: 10.1175/1520-0485(1989)
 503 019<0551:otrooe>2.0.co;2
- 504 Beyer, K., Goldstein, J., Ramakrishnan, R., & Shaft, U. (1998). When is “near-
 505 est neighbor” meaningful? In *Lecture notes in computer science (including*
 506 *subseries lecture notes in artificial intelligence and lecture notes in bioinfor-*
 507 *matics)*. doi: 10.1007/3-540-49257-7{_}15
- 508 Buizer, J. L., Foster, J., & Lund, D. (2000). Global impacts and regional actions:
 509 Preparing for the 1997-98 El Niño. *Bulletin of the American Meteorological So-*
 510 *ciety*. doi: 10.1175/1520-0477(2000)081<2121:GIARAP>2.3.CO;2
- 511 Dijkstra, H. A., Petersik, P., Hernández-García, E., & López, C. (2019, 10). The Ap-
 512 plication of Machine Learning Techniques to Improve El Niño Prediction Skill.

- 513 *Frontiers in Physics*, 7. Retrieved from [https://www.frontiersin.org/](https://www.frontiersin.org/article/10.3389/fphy.2019.00153/full)
514 [article/10.3389/fphy.2019.00153/full](https://www.frontiersin.org/article/10.3389/fphy.2019.00153/full) doi: 10.3389/fphy.2019.00153
- 515 Emery, W. J. (2008). Water Types and Water Masses. In *Encyclopedia of ocean sci-*
516 *ences: Second edition*. doi: 10.1016/B978-012374473-9.00108-9
- 517 Iizumi, T., Luo, J. J., Challinor, A. J., Sakurai, G., Yokozawa, M., Sakuma, H., ...
518 Yamagata, T. (2014). Impacts of El Niño Southern Oscillation on the global
519 yields of major crops. *Nature Communications*. doi: 10.1038/ncomms4712
- 520 IPCC. (2019). IPCC Special Report on the Ocean and Cryosphere in a Changing
521 Climate. In *Ipcc summary for policymakers*. doi: <https://www.ipcc.ch/report/srocc/>
522
- 523 Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern Recognition*
524 *Letters*. doi: 10.1016/j.patrec.2009.09.011
- 525 Johnson, G. C., & Birnbaum, A. N. (2016). Equatorial Pacific Thermostat re-
526 sponse to El Niño. *Journal of Geophysical Research: Oceans*. doi: 10.1002/
527 2016JC012304
- 528 Johnson, G. C., & Birnbaum, A. N. (2017). As El Niño builds, Pacific Warm Pool
529 expands, ocean gains more heat. *Geophysical Research Letters*. doi: 10.1002/
530 2016GL071767
- 531 Jones, D. C., Holt, H. J., Meijers, A. J., & Shuckburgh, E. (2019). Unsupervised
532 Clustering of Southern Ocean Argo Float Temperature Profiles. *Journal of*
533 *Geophysical Research: Oceans*. doi: 10.1029/2018JC014629
- 534 Karpatne, A., Ebert-Uphoff, I., Ravela, S., Babaie, H. A., & Kumar, V. (2019).
535 Machine Learning for the Geosciences: Challenges and Opportunities.
536 *IEEE Transactions on Knowledge and Data Engineering*. doi: 10.1109/
537 TKDE.2018.2861006
- 538 Karsten, R. H., & Marshall, J. (2002). Constructing the residual circulation of the
539 ACC from observations. *Journal of Physical Oceanography*. doi: 10.1175/1520-
540 -0485(2002)032<3315:CTRCOT>2.0.CO;2
- 541 Kim, K. Y., & Kim, Y. Y. (2002). Mechanism of Kelvin and Rossby waves during
542 ENSO events. *Meteorology and Atmospheric Physics*. doi: 10.1007/s00703-002-
543 -0547-9
- 544 Marshall, J., & Plumb, R. A. (2008). *Atmosphere, Ocean, and Climate Dynamics*.
545 doi: 10.1017/CBO9781107415324.004

- 546 Maze, G., Mercier, H., & Cabanes, C. (2017). Profile Classification Models. *Mercator*
547 *Ocean Journal*.
- 548 Maze, G., Mercier, H., Fablet, R., Tandeo, P., Lopez Radcenco, M., Lenca, P., ...
549 Le Goff, C. (2017). Coherent heat patterns revealed by unsupervised classification
550 of Argo temperature profiles in the North Atlantic Ocean. *Progress in*
551 *Oceanography*. doi: 10.1016/j.pocean.2016.12.008
- 552 Meinen, C. S., & McPhaden, M. J. (2000). Observations of warm water vol-
553 ume changes in the equatorial Pacific and their relationship to El Nino and
554 La Nina. *Journal of Climate*. doi: 10.1175/1520-0442(2000)013<3551:
555 OOWWVC>2.0.CO;2
- 556 NOAA National Centers for Environmental Prediction. (2020). *Cold & Warm*
557 *Episodes by Season*. Retrieved from [https://origin.cpc.ncep.noaa.gov/
558 products/analysis_monitoring/ensostuff/ONI_v5.php](https://origin.cpc.ncep.noaa.gov/products/analysis_monitoring/ensostuff/ONI_v5.php)
- 559 Pedregosa, F., Michel, V., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., ...
560 Duchesnay, E. (2011). *Scikit-learn: Machine Learning in Python* (Vol. 12;
561 Tech. Rep.). Retrieved from <http://scikit-learn.sourceforge.net>.
- 562 Rand, W. M. (1971). Objective criteria for the evaluation of clustering methods.
563 *Journal of the American Statistical Association*. doi: 10.1080/01621459.1971
564 .10482356
- 565 Rasmusson, E. M., & Carpenter, T. H. (1982). Variations in tropical sea
566 surface temperature and surface wind fields associated with the South-
567 ern Oscillation/ El Nino (Pacific) . *Monthly Weather Review*. doi:
568 10.1175/1520-0493(1982)110<0354:VITSST>2.0.CO;2
- 569 Riser, S. C., Freeland, H. J., Roemmich, D., Wijffels, S., Troisi, A., Belbéoch, M., ...
570 Jayne, S. R. (2016). *Fifteen years of ocean observations with the global Argo*
571 *array*. doi: 10.1038/nclimate2872
- 572 Roemmich, D., & Gilson, J. (2011). The global ocean imprint of ENSO. *Geophysical*
573 *Research Letters*. doi: 10.1029/2011GL047992
- 574 Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and vali-
575 dation of cluster analysis. *Journal of Computational and Applied Mathematics*.
576 doi: 10.1016/0377-0427(87)90125-7
- 577 Shlens, J. (2003). A tutorial on principal component analysis: derivation, discussion
578 and singular value decomposition. *Online Note [httpwww snl salk edushlenspub-](http://www.snl.salk.edu/shlenspub)*

579 *notespca pdf*. doi: 10.1.1.115.3503

580 Timmermann, A., An, S. I., Kug, J. S., Jin, F. F., Cai, W., Capotondi, A.,

581 ... Zhang, X. (2018). *El Niño–Southern Oscillation complexity*. doi:

582 10.1038/s41586-018-0252-6

583 Wang, B., Wu, R., & Lukas, R. (2000). Annual adjustment of the thermocline in
584 the tropical Pacific Ocean. *Journal of Climate*. doi: 10.1175/1520-0442(2000)

585 013(0596:AAOTTI)2.0.CO;2

586 Wyrтки, K. (1975). El Niño—The Dynamic Response of the Equatorial Pacific
587 Ocean to Atmospheric Forcing. *Journal of Physical Oceanography*. doi:

588 10.1175/1520-0485(1975)005(0572:entdro)2.0.co;2

589 Xu, D., & Tian, Y. (2015). A Comprehensive Survey of Clustering Algorithms. *An-*
590 *nals of Data Science*. doi: 10.1007/s40745-015-0040-1

591 Yang, H., & Wang, F. (2009). Revisiting the thermocline depth in the equatorial Pa-
592 cific. *Journal of Climate*. doi: 10.1175/2009JCLI2836.1

593 Yeh, S. W., Kug, J. S., Dewitte, B., Kwon, M. H., Kirtman, B. P., & Jin, F. F.

594 (2009). El Niño in a changing climate. *Nature*. doi: 10.1038/nature08316

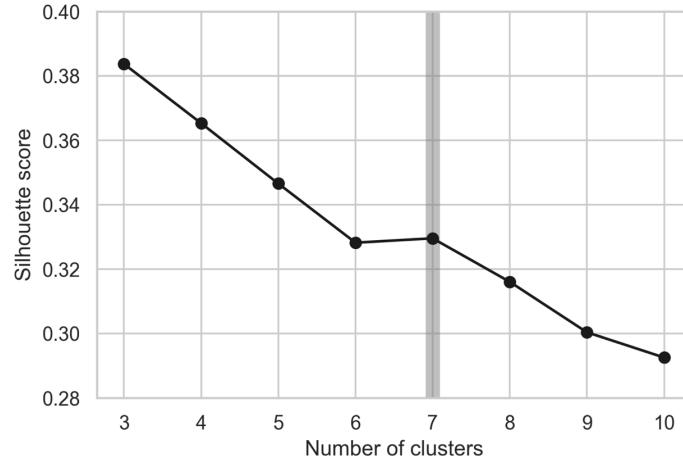


Figure 1. Silhouette score as a function of number of clusters, k , from 3 to 10 calculated following equation 4. A local maximum (highlighted in gray) is observed at $k = 7$.

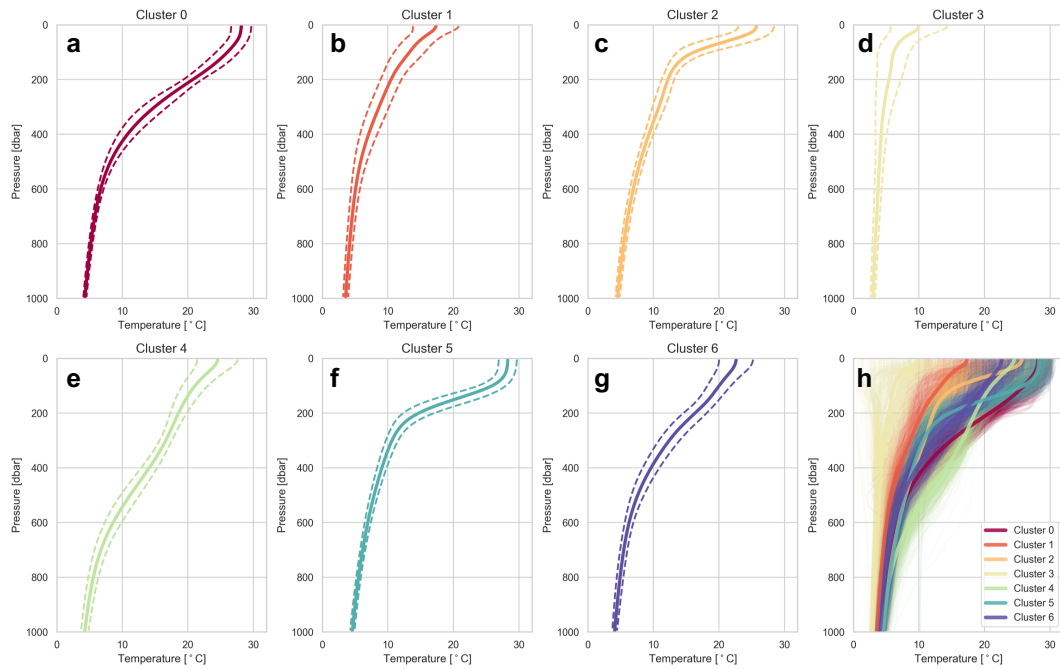


Figure 2. (a-g) For each cluster, the mean temperature profile (solid line) and \pm one standard deviation of temperature (dashed line) is plotted. (h) Overlay of a random subset of profiles from each cluster, with thicker lines indicating the mean temperature profile in each cluster, colored by cluster assignment.

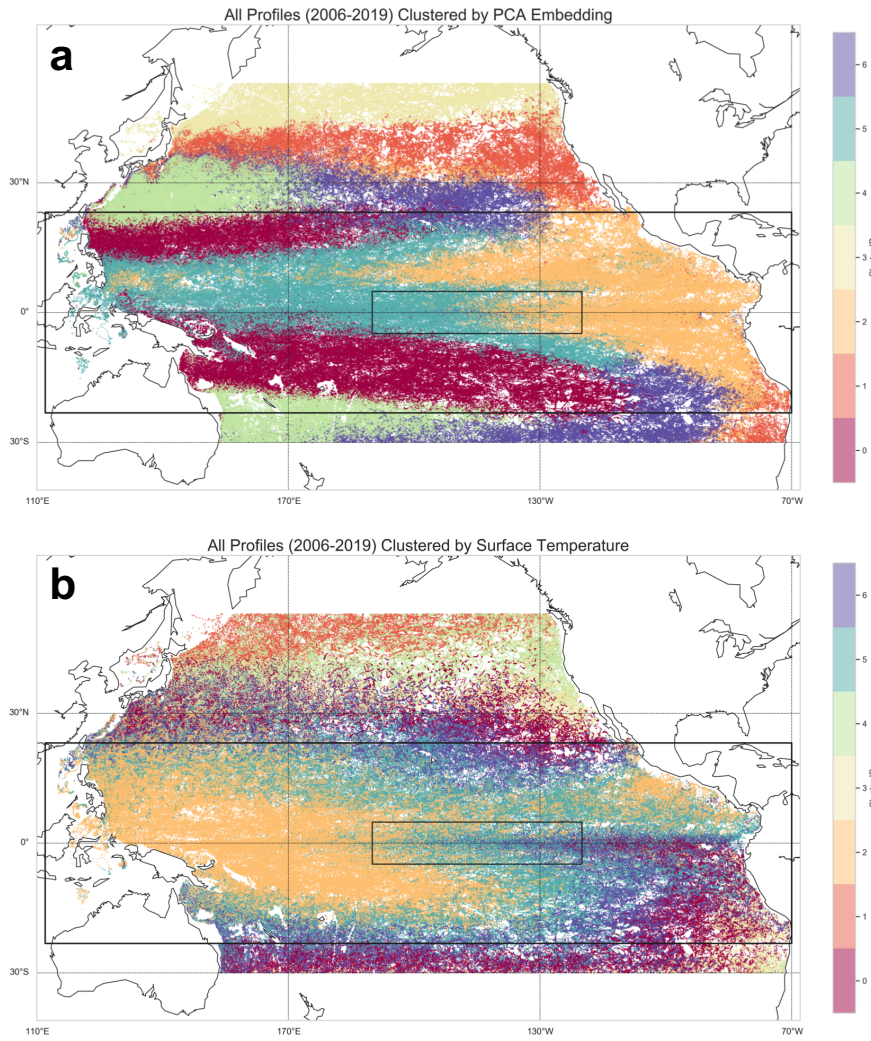


Figure 3. The spatial distribution of Argo measurements in the Pacific, colored by cluster assignment. Cluster IDs are randomly set by the clustering algorithm initialization, therefore ID magnitudes are arbitrary. The large black box corresponds to the tropical zone ($\pm 23.4^\circ$ latitude), and the smaller inner box corresponds to the Niño 3.4 region. (a) Measurements grouped by PCA embedding of full temperature profile, used for subsequent analysis. (b) Measurements grouped by sea surface temperature (uppermost profile measurement only).

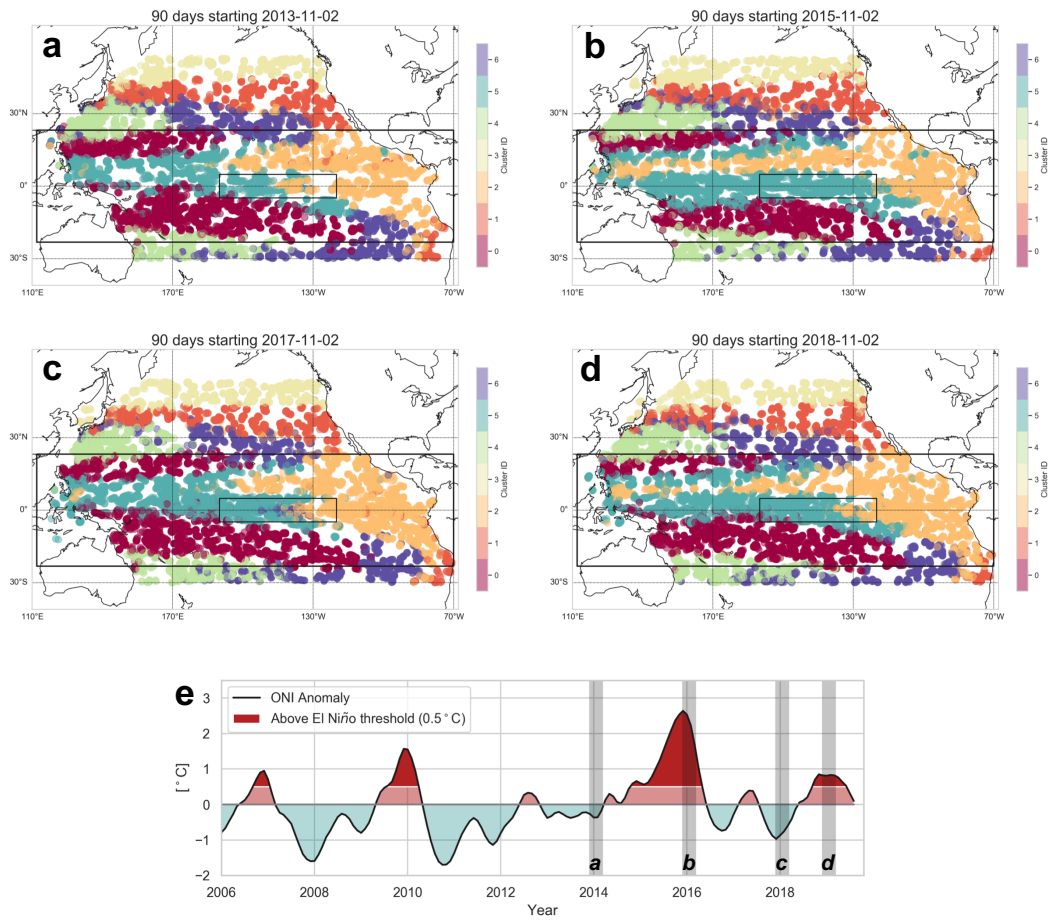


Figure 4. (a-d) Measurements within each three-month period indicated colored by cluster assignment (see supplementary video for cluster assignments over all time). (e) The ONI anomaly from 2006 to 2019 indicating several El Niño events. Vertical gray shaded bars correspond to time periods visualized in upper plots.

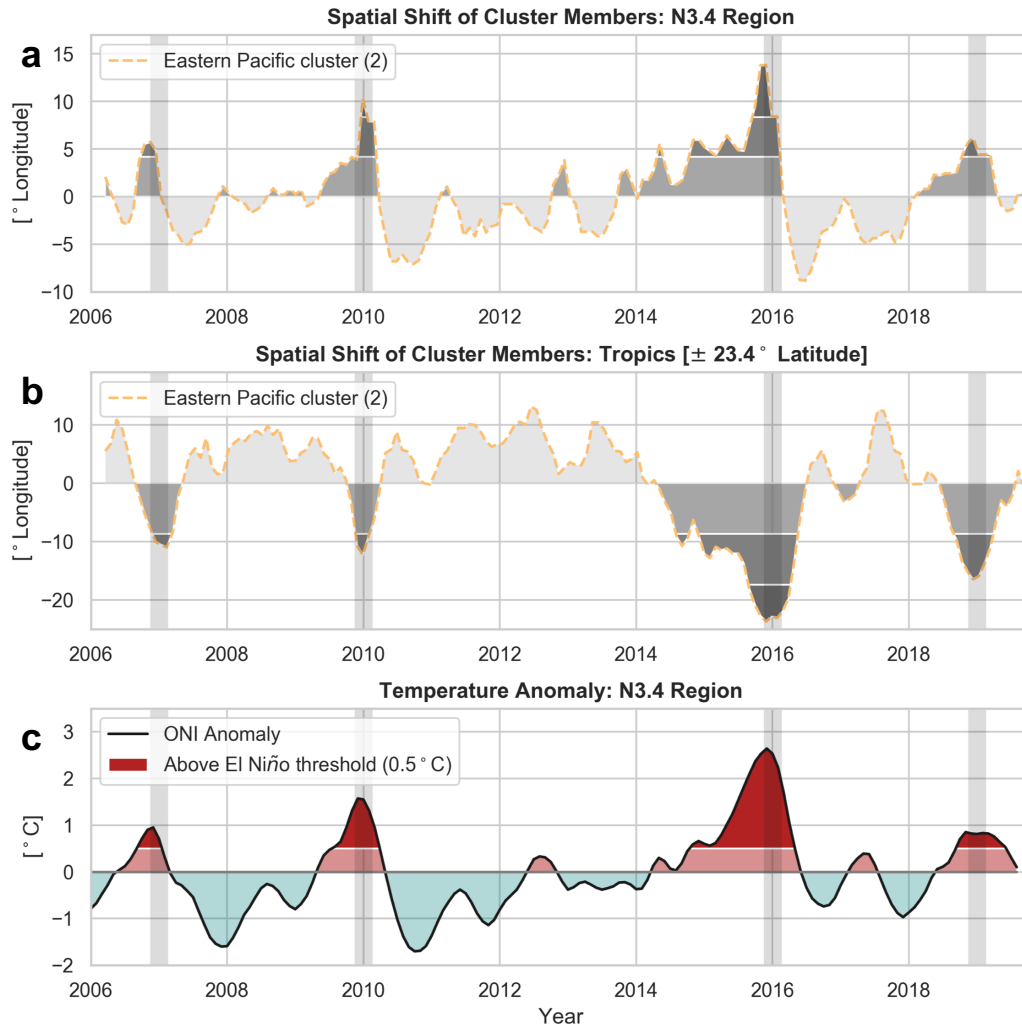


Figure 5. (a) The longitudinal anomaly of the eastern cluster members within the Niño 3.4 region. (b) The longitudinal anomaly of the Eastern cluster members over the entire tropics. White lines and gray shading correspond to standard deviations from the mean. Vertical gray bars on all plots correspond to an El Niño event occurring. (c) ONI during the same period, dark red region corresponds to events above 0.5°C threshold.