



Filtering method based on cluster analysis to avoid salinity drifts and recover Argo data in less time

Emmanuel Romero¹, Leonardo Tenorio-Fernandez², Iliana Castro¹, and Marco Castro¹

¹Tecnológico Nacional de México/Instituto Tecnológico de La Paz, La Paz, México

²CONACyT-Instituto Politécnico Nacional-Centro Interdisciplinario de Ciencias Marinas, La Paz, México

Correspondence: Leonardo Tenorio-Fernandez (ltenoriof@ipn.mx)

Received: 11 March 2021 – Discussion started: 6 April 2021

Revised: 17 July 2021 – Accepted: 27 July 2021 – Published: 17 September 2021

Abstract. Currently there is a huge amount of freely available hydrographic data, and it is increasingly important to have easy access to it and to be provided with as much information as possible. Argo is a global collection of around 4000 active autonomous hydrographic profilers. Argo data go through two quality processes, real time and delayed mode. This work shows a methodology to filter profiles within a given polygon using the odd–even algorithm; this allows analysis of a study area, regardless of size, shape or location. The aim is to offer two filtering methods and to discard only the real-time quality control data that present salinity drifts. This takes advantage of the largest possible amount of valid data within a given polygon. In the study area selected as an example, it was possible to recover around 80 % in the case of the first filter that uses cluster analysis and 30 % in the case of the second, which discards profilers with salinity drifts, of the total real-time quality control data that are usually discarded by the users due to problems such as salinity drifts. This allows users to use any of the filters or a combination of both to have a greater amount of data within the study area of their interest in a matter of minutes, rather than waiting for the delayed-mode quality control that takes up to 12 months to be completed. This methodology has been tested for its replicability in five selected areas around the world and has obtained good results.

1 Introduction

Autonomous oceanographic instruments have become very important tools in observational oceanography. Hydrographic autonomous profilers (HAPs) are tools that reduce the costs of in situ oceanographic observations, obtaining a large number of hydrographic profiles in time and space, at a lower cost compared to those carried out on oceanographic cruises. An example of these HAPs is those belonging to the Argo program, in which each measured profile is processed by its Data Assembly Center (DAC) in a quality control system, before being published (Argo Data Management Team, 2019).

HAPs have the ability to continuously measure hydrographic parameters in the water column. Since the beginning of the program and up to the present, there are data records collected from around 15 300 core HAPs and around 1300 biogeochemical HAPs belonging to the global Argo group in the world oceans, which have measured temperature, salinity and biogeochemical parameters in most cases from 2000 m depth to the sea surface or vice versa, from which around 4000 are currently active (Argo, 2020a). However, around 75 % of the total profiles have completed the quality control process, and therefore it is considered that the rest may or may not be of such good quality. In areas with a low concentration of profiles, this percentage is more significant, and it is important to obtain as much data as possible to support scientific research.

The data of each HAP have to be validated and processed by a quality control system, before being used or published. The Argo quality control system consists of two stages, real-time quality control (RTQC) and delayed-mode quality control (DMQC). The tests performed by the RTQC are automated and limited, due to the requirement to be available within the first 24 h after transmission. These data are free of serious errors in each of their variables (e.g. impossible data in dates and coordinates) and must be within the global and regional ranges. In the case of having adjusted parameters available, these are placed in the same variables, but named with the suffix “_ADJUSTED”; in this way the data are preserved without adjustments in the variables without this suffix. The second quality control process is the DMQC. The data adjusted by this quality control replace the data adjusted by the RTQC. Since, during this process, the data are subjected to detailed scrutiny by oceanographic experts, DMQC data can take a year to be published (Wong et al., 2021). Normally, due to the problems presented by the RTQC data, such as the salinity drifts presented in this work, users of Argo program data are advised to only use DMQC data for scientific analysis, or to perform it by themselves; this quality control is explained in the manuals, and for this reason many users decide not to use the RTQC data.

The objective of this work is to present a methodology based on cluster analysis to admit the data in RTQC that conforms to the same hydrography patterns as the DMQC data and thus increase the amount of data available for scientific research, avoiding the complete discard of the RTQC data. To carry out this methodology, first, the data must be delimited by a polygon, the one that represents the study area of interest. Using a point-in-polygon (PIP) algorithm the profiles that were measured within the study area of interest are determined. In addition, a web application was developed to show the results of the application of this methodology and the usefulness that it can have if it were integrated into the HTTP data access platforms, such as statistics and graphs of study areas defined by the user.

2 Data collection and methods

To achieve the objectives of this research in any study area given by a polygon, irregular or not, and since the selection of the data can be of interest both at a global and regional level, the geographical coordinates of the profiles stored on the servers of Argo were used as points for the PIP problem and thus determine if they were measured within the study area. To solve it, the even–odd algorithm (Foley and Hughes, 1990) was used, and once the profiles within the polygon are obtained, the profile data are downloaded.

For the purposes of testing the methods of this work, a study area was selected (Fig. 1), which is located between 25 and 19° N and 113 and 105° W. In this area it is known that there are current interactions between the tropical branch of

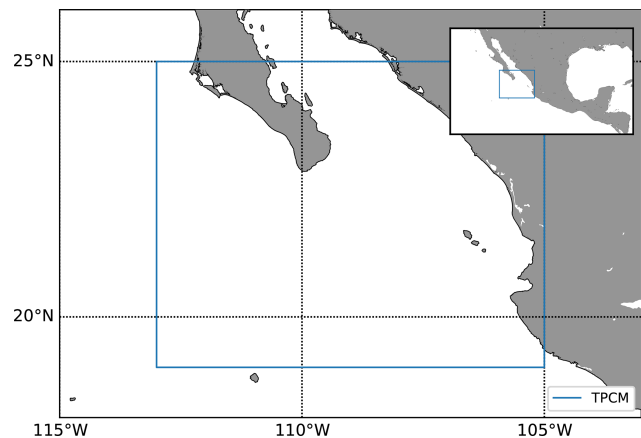


Figure 1. Study area. The upper right corner shows the location of the study area composed of parts of the California Current System, the Gulf of California, the transition area and the tropical Pacific off central Mexico (TPCM), shown in the foreground.

the California Current, the Gulf of California current and to the north the Mexican Coastal Current (Lavín and Marinone, 2003; Lavín et al., 2009; Godínez et al., 2010; Portela et al., 2016). These interactions produce intense mesoscale activity and a high complexity in the circulation (Kessler, 2006). Two mesoscale structures such as cyclonic and anticyclonic eddies interact and play an important role in circulation (Zamudio et al., 2001, 2007; Lavín et al., 2006; Pantoja et al., 2012). As well as this, this area is part of the minimum oxygen zone (Fiedler and Talley, 2006; Stramma et al., 2008). The study area encompasses parts of the California Current System, the Gulf of California, the transition area and the tropical Pacific off central Mexico (Portela et al., 2016); hereinafter this study area will be named TPCM.

One of the great benefits of using a PIP algorithm to filter locations is that it can be used with data from other georeferenced databases. To demonstrate this, tests were carried out with the World Ocean Atlas 2018 (WOA18) database, which provides quality-controlled data to calculate the climatology of temperature, salinity, dissolved oxygen and dissolved inorganic nutrients derived from profiling floats, OSD (ocean station data), CTD (conductivity, temperature and depth) and many other contents in the NCEI World Ocean Database 2018 (WOD18). The monthly data of temperature (Locarnini et al., 2018) and salinity (Zweng et al., 2018) of the statistical mean of each quarter degree ($1/4^\circ$) from 2005 to 2017 were downloaded, and the PIP algorithm described in this work was applied to the polygon that delimits the TPCM. Conversions from in situ temperature to conservative temperature and practical salinity to absolute salinity were carried out according to the Thermodynamic Equation of SeaWater 2010 (TEOS-10). The monthly T – S (temperature and salinity) diagrams of the WOA18 data and the Argo DMQC data were compared to corroborate that they are located in the

same water masses and to review the quality of the DMQC data in the area.

The data measured within the TPCM were statistically analyzed, and it was found that there are few profiles within the area and that around only 30 % of the total data have been evaluated by the RTQC. The Argo manual (Argo Data Management Team, 2019) indicates that the quality control flags establish how good or bad the data are, with 1 being good data and, as the value increases, the quality deteriorates. These flags are determined in the RTQC, and when the DMQC arrives it creates flags for the adjusted data. In the event that data have been adjusted in the RTQC, they are replaced, because during the RTQC only gross errors can be identified by the automated procedures, and to detect more subtle issues with sensor drift these can only be carried out in DMQC (Wong et al., 2021). Tests were performed by graphing the T – S diagrams using these flags and adding the density isoline and the water masses according to Portela et al. (2016). The RTQC data of all the flags were used, and all of them showed salinity drifts, including the marked data with flag 1 as seen in supplementary material A.1 (Romero et al., 2021b), so it is not feasible to use these indicators to filter the data in RTQC. To increase the amount of available data, cluster analysis was applied to the data, since two groups of data can be visually located in the T – S diagrams: those that form the same patterns as those of the DMQC and those that do not. This analysis groups a set of objects in such a way that the characteristics of the objects of the same group are more similar to each other than to the other groups (Everitt et al., 2011). In this case, the aim is to separate the RTQC data into groups, a group that contains data with characteristics similar to DMQC data and other groups with salinity drift problems.

To perform the cluster analysis, the unsupervised k -means classification algorithm was chosen, which groups the data into k groups, minimizing the distance between the data and the centroid of its group (Hartigan and Wong, 1979). The algorithm starts by setting the k centroids in the data space, regardless of where the data were obtained, and assigning the data to their closest centroid. Then, it updates the position of the centroid of each group, calculating the position of the average of the data belonging to each group, and the data are reassigned to their closest centroid. This process is repeated until the centroids do not change position. An algorithm based on distances was selected because it seeks to obtain only the RTQC data closest to the DMQC data.

Since it is necessary to indicate the number of k centroids when we use k -means, a manual enumeration of the groups to be searched is required. To automate this process and avoid the user having to indicate the exact number of centroids needed to retrieve RTQC data for each month and for each study area chosen, Algorithm 1 was programmed.

Algorithm 1 receives the adjusted data from the DMQC and the RTQC; in the case of profiles that have not been adjusted, the data are received without adjustment. The algo-

ALGORITHM 1
 RTQC data filtering

```

dataset ← FilterByMonth(dataset)
for i ← 0 to 11 do
  for j ← 0 to 10 do
    data ← GetDataWithDepthHigherThan(dataset[i], depth
      ← 1500)
    mid_ranges ← GetMidRangeOfDMQCandRTQC(data)
    groups ← kmeans(data, k ← 2, init ← mid_ranges)
    if groups[0] have DMQC data and groups[1] do not then
      dataset[i] ← MatchDataByProfilerAndProfile(dataset[i],
        groups[0])
    else
      break
    end if
  end for
end for
return dataset

```

rithm separates these data by month in an array and iterates them. Within each iteration, it calculates the salinity mid-range of each quality control and divides the data measured at depths greater than that specified by the user, 1500 m by default, into two groups (using the mid-ranges as the starting position of the centroids), up to a maximum of 10 iterations, each time verifying if there are DMQC data in both groups. If so, the algorithm stops and returns the data without grouping them; on the contrary, if only a group contains the data in DMQC, it associates the data of that group with the data at depths less than that specified by the user, taking into consideration the month, the profiler code and the profile number, and replaces the group data with the associated data. The mid-ranges are used as the initial position of the centroids to prevent them from being generated randomly. The procedure described above is the first filter of the RTQC data; in each iteration the algorithm discarded the groups that presented salinity drifts and kept only the group where the DMQC data were found, and thus, when the execution of the algorithm ends, RTQC data within the group with DMQC data are considered to contain no salinity drifts. To increase the reliability of the filtering, a second filter was created. In the second filter, the algorithm stores in memory the profilers that presented salinity drifts during the execution of the Algorithm 1. Thus, the second filter not only discards the profiles with drift problems, it also directly discards all the profiles of the profilers that they have at least one profile with salinity drifts.

A library was developed that contains all the procedures described in this work. Using it, as an extra example, five study areas were delimited with different extensions, locations, profile densities and hydrographic characteristics. The first area was the Alboran Sea, which was selected because the data were measured in shallow water (0 to 1200 m). The Antarctic area was selected because of its high latitude (cold water). The third area was the Bermuda Triangle, which was

selected because it is located in the Atlantic transition between the tropical and subtropical area. The fourth is the tropical zone of the Pacific that surrounds an archipelago of the central Pacific, and the last one is in the tropical sea of Indonesia. All data from these areas were downloaded from the snapshot of December 2020 (Argo, 2020b) and evaluated by Algorithm 1.

To test the above methods in a more extensive and irregular polygon area, a web application was developed. The study area for this web application was delimited using the Exclusive Economic Zone (EEZ) of Mexico as an example, and the geographical location of the profiles from around the world are filtered by the PIP algorithm, to automatically download the data every 24 h within this irregular polygon through the IFREMER synchronization service (available at <http://www.argodatamgt.org/Access-to-data>, last access: November 2019).

Every time that new data from HAPs are downloaded, they go through a processing phase, in which they are cleaned and transformed to be integrated into the web application. For example, the variables of temperature and salinity are converted to conservative temperature and absolute salinity, as defined by TEOS-10; the current description of the properties of seawater defines it. Afterwards, graphs and useful files are generated to show information about the HAPs and their profile data.

The web application was developed on a satellite map, to which tools were added for data management and visualization, such as drawing irregular polygons to define study areas within the main polygon, filtering data to display statistical and graphical data according to the selected filter and trajectory tracing, among others. Also, RTQC data filtering was implemented in the web application. The same irregular polygons with which statistical data are obtained can be used to indicate a study area in which as much data as possible is obtained without salinity drifts.

3 Results

The data used to obtain the following results were downloaded from the Coriolis GDAC FTP server in 2019 and the “Profile directory file of the Argo Global Data Assembly Center” file was used as input for the chosen PIP algorithm, which filtered the measured profiles inside the polygon correctly. Figure 2 shows the result of the T – S diagram comparison between the DMQC data and the WOA18 data. The DMQC and WOA18 data are located in the same water masses, and the data are spliced at depths greater than 1500 m, which validates that the DMQC data follow the same patterns as the data from other international databases. According to Portela et al. (2016), this region is made up of the California Current Water (CCW), Tropical Surface Water (TSW), Gulf of California Water (GCW), Subtropical Sub-surface (SS) and the Pacific Intermediate Water (PIW).

On the contrary, the data in RTQC with the best quality flag present drifts in salinity. The RTQC and DMQC data were plotted in the T – S diagrams together per month of the TPCM, and some of the data in RTQC were the cause of salinity drifts in almost all the months (Fig. 3).

In Fig. 3 it is clear that the salinity drifts in the RTQC data are important; however, it is also shown that certain data follow the structure (shape) of the DMQC data. To avoid discarding all RTQC data, the use of cluster analysis is proposed. By applying cluster analysis to all data in RTQC with the k -means algorithm and with different values in k , the resulting groups mix data that show salinity drifts with data that follow the same patterns as the DMQC data at 1500 m. This is because, at depths less than 1500 m, salinity data are more dispersed than at greater depths.

Taking into consideration that at depths greater than 1500 m, the variations in salinity and temperatures are imperceptible in this study area, the cluster analysis was performed with the salinity data measured at depths greater than 1500 m. The resulting groups are shown in Fig. 4a and b, and in the figure it can be observed that one of the resulting groups contains the data that follow the same patterns as the DMQC data, and the rest of the groups contain data with salinity drifts. Therefore the next step was to associate the data of these groups with the rest of the data, taking into consideration the profiler code and the profile number and thus obtaining complete groups (Fig. 4c and d).

Figure 4 shows how the groups are separated with the chosen algorithm. In the months of January and November, DMQC data are displayed as yellow dots and the orange groups contain the RTQC data that follow the patterns of the data in DMQC. The blue, green and red groups contain the data showing salinity drifts.

To avoid indicating the number of k centroids manually, Algorithm 1 was developed. Figure 5 shows the first three iterations of the month of November as an example. In Fig. 5a and b blue data represent the group that contains DMQC data and the orange color group represents the group of the RTQC data. The data contained in the orange groups are discarded by the algorithm. Figure 5c is the third iteration, and both groups contain data in DMQC, because the data are so close to each other that the k -means algorithm (which is based on distances to separate the groups) divides the DMQC data into two different groups, so that in this iteration the algorithm stops.

The results of the first filtering of the proposed algorithm are shown in Fig. 6a. The filtered data from the RTQC show the same patterns as the DMQC data, except for the months of July, August and September. In July and August, the salinity drifts are found at depths less than 1500 m, while in September, the drifts present values very close to the DMQC data and this prevents the algorithm from being able to separate them. This filter allows a greater amount of admitted RTQC data to be obtained, but as seen in the figure, it still shows salinity drifts in some cases. For this reason, the sec-

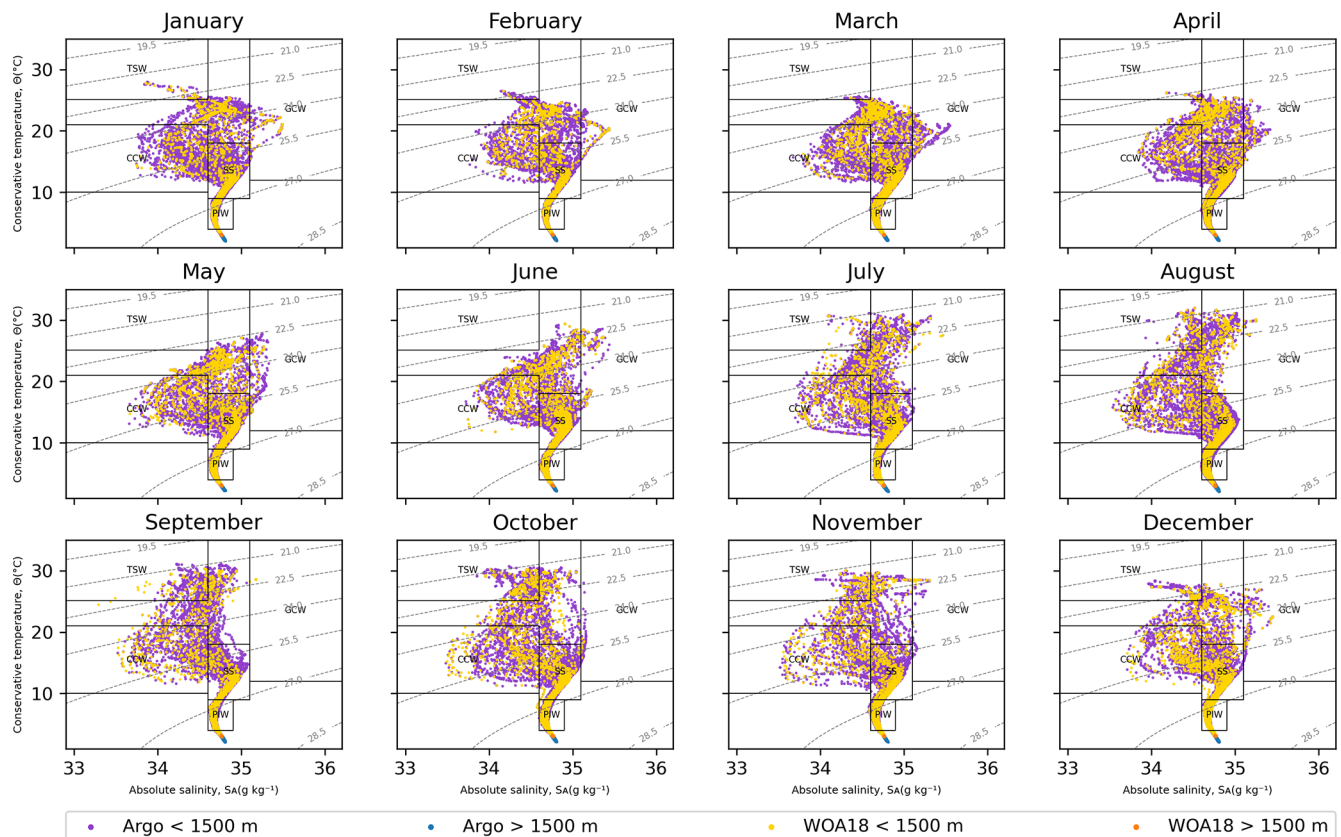


Figure 2. Monthly comparison of T – S diagrams of data from DMQC and WOA18. The black boxes delimit the limits of the water masses in the region and the gray isolines the density (kg m^{-3}).

ond filter was incorporated. Figure 6b shows the results of it, since it considers those profilers that have presented salinity drifts and removes their profiles completely, a significant reduction in admitted data from the RTQC is observed, but these no longer show salinity drifts.

Table 1 shows the total measurements (Meas.) made in the TPCM area and the measurements filtered by the aforementioned algorithms.

The total usable data in the TPCM due to the first and second filters represent $\sim 95\%$ and $\sim 80\%$ of the data, compared to the $\sim 70\%$ that would be obtained by automatically discarding the data in RTQC. By presenting this option to the user and filtering the data from the RTQC, instead of discarding $\sim 30\%$ of the total, only $\sim 5\%$ would be discarded in the case of the first filter and $\sim 20\%$ in the case of the second, which would mean a considerable increase in the data available for use. After all, the admitted data present similar characteristics to the data that were already evaluated with the DMQC. They have a high probability of not needing adjustments and therefore could be used in research before waiting for the DMQC to be applied to them.

Despite the fact that in the first filter some months were not filtered in the desired way in the study area, the user may simply not use the data from those months or use the second

filter if the user wishes to use only the most reliable data. Also, the possibility of using a combination of both filters is not ruled out – if the user uses the months of the first filter that no longer present salinity drifts and uses the data of the second filter in which they present drifts, the largest possible amount of admissible data would be used in any study area. The results of the algorithm will change depending on the extension and the hydrographic characteristics of the study area that the user selects. Selecting which filter to use or whether to make a combination of them, as well as deciding whether to use the default depth or use a more suitable one for the study area, is the responsibility of the user, and it is recommended to have knowledge of the study area.

A library for Python 3.7 named *cluster_qc* was developed alongside this work. It contains all the procedures described here and is available under the Creative Commons Attribution 4.0 International License (latest package version is v1.0.2; Romero et al., 2021a). Using this library, five study areas were delimited with different extension, location, profile density and hydrographic characteristics, and the data were downloaded from the snapshot of December 2020 (Argo, 2020b) and evaluated by Algorithm 1. The results are shown in Table 2, and the figures of these results in supplementary material A.2 (Romero et al., 2021b).

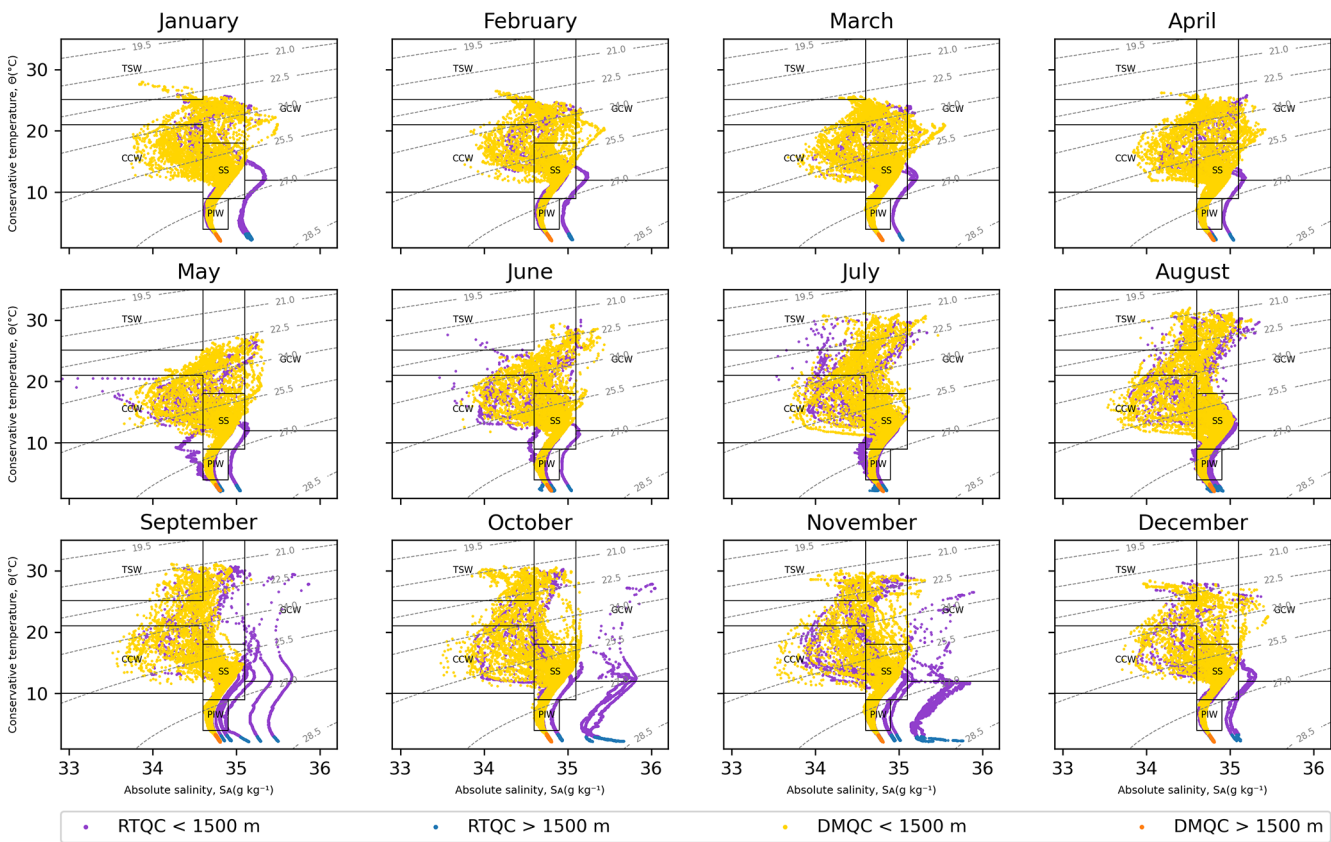


Figure 3. Monthly comparison of T – S diagrams of data from RTQC and DMQC. The black boxes delimit the limits of the water masses in the region and the gray isolines the density (kg m^{-3}).

Table 1. Percentages of DMQC and RQTC data admitted and discarded normally and by the two proposed filters.

Data	Without filter		First filter		Second filter	
	Meas.	%	Meas.	%	Meas.	%
DMQC	594 385	69.96 %	594 385	69.96 %	594 385	69.96 %
Admitted RTQC	0	0.00 %	209 392	24.64 %	82 196	9.67 %
Discarded RTQC	255 184	30.03 %	45 792	5.39 %	172 988	20.36 %
Total	849 569	100.00 %	849 569	100.00 %	849 569	100.00 %

In the results of the Alboran Sea, the westernmost part of the Mediterranean Sea, there are no data deeper than 1500 m or salinity drifts, so the algorithm directly returns the data set without modification. The algorithm receives the depth of 1500 m by default, and sending a lower depth could eliminate salinity variations if there were any. In the case of Antarctica, we found that the months of February and April contain salinity drifts, which could not be completely eliminated with the first filter. For this case, it is recommended to use the RTQC data supported by the second filter. On the other hand, in the Bermuda Triangle, salinity drifts are shown in the months of June to October, in addition to atypical values in the rest of the months. The first filter already eliminates

salinity drifts, so in this case it is recommended to use this filter and eliminate outliers. In the fourth study area, which surrounds a central Pacific archipelago, there are many outliers in all the months; however, the first filter managed to rule out the salinity drifts present in the months of September to December. In this case it is recommended to reduce the study area into smaller areas to apply the filters and treat the outliers separately. Finally, the large study area located next to Indonesia shows salinity drifts in the months of March and July to December. The first filter was able to filter the salinity drifts except for the month of December, because the deviations are above 1500 m. In this case it is recommended to

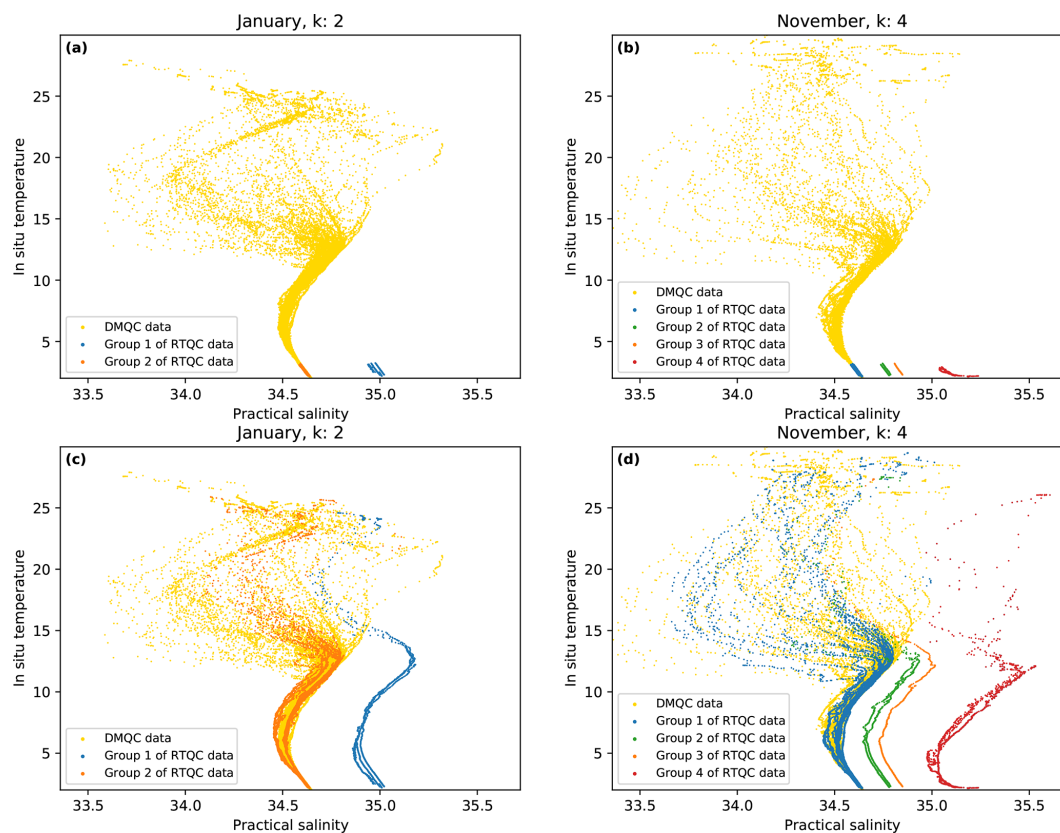


Figure 4. Cluster analysis results. Panels (a) and (b) show the groups formed with the RTQC data measured at depths greater than 1500 m. Panels (c) and (d) show these same grouped data but matched data with the rest of their profile data.

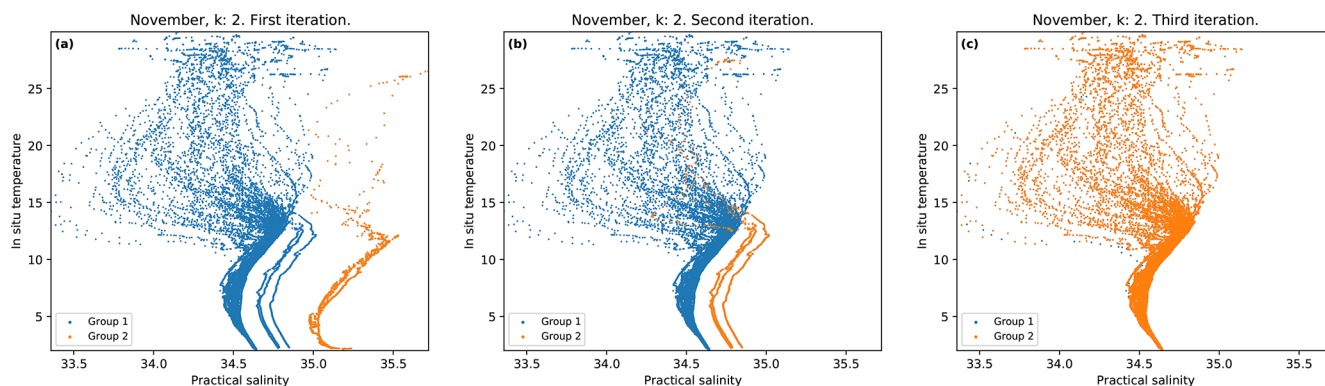


Figure 5. First three iterations of the proposed algorithm using the data for the month of November. Panels (a), (b) and (c) are the first, second and third interactions.

use the data from the first filter for the months of January to November, or use only the months with no outliers.

3.1 Web application

The web application got interesting results and can be accessed through the *cluster_qc* library repository. In Fig. 7, it is observed that the PIP algorithm filters the profiles that were made within the EEZ of Mexico correctly. The blue

line represents the given polygon, and the locations of the filtered profiles inside and outside the polygon are represented by dots in red and black respectively.

Once the data have been downloaded and transformed, statistical data specific to the EEZ of Mexico can be obtained, such as the number of profilers within the polygon, the number of profiles or profilers per year, or the DACs to which these profilers belong, among others. Table 3 shows the profilers that have carried out measurements within the poly-

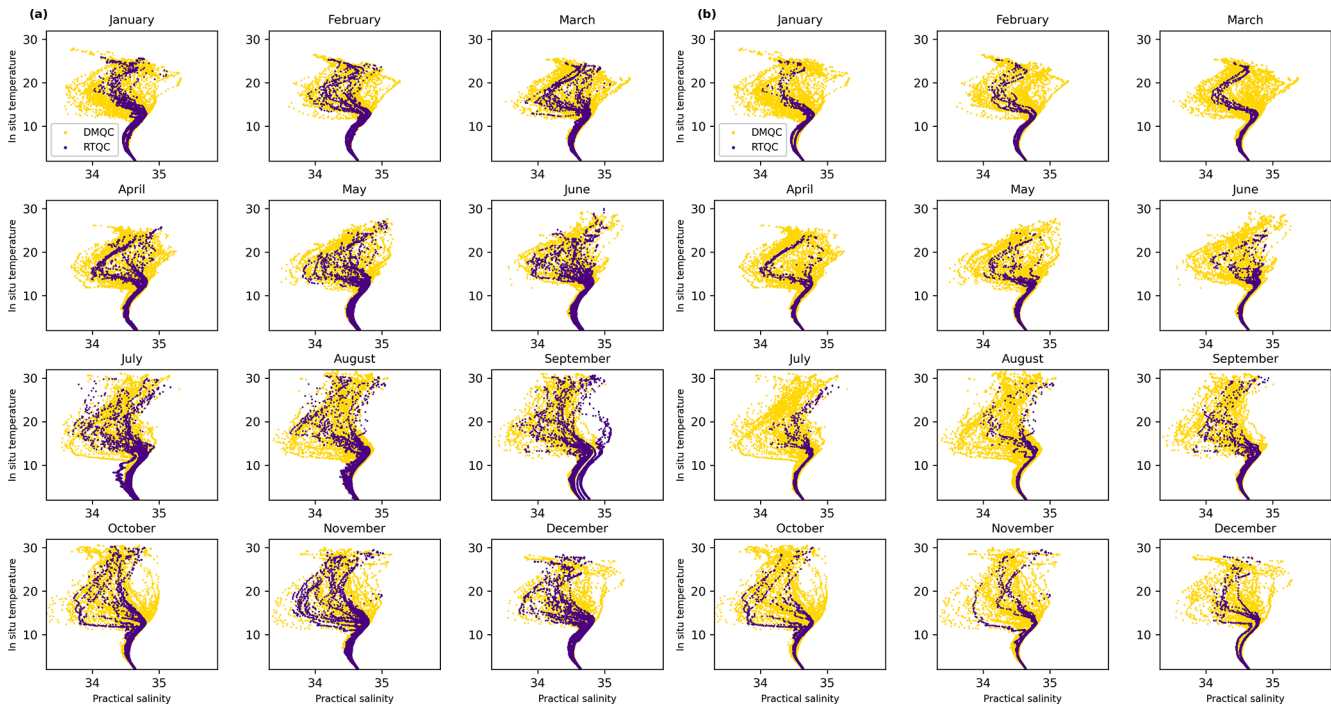


Figure 6. Monthly comparison of the T – S diagrams of RTQC and DMQC. (a) First filtering of RTQC. (b) Second filtering of RTQC.

Table 2. Results of Algorithm 1 in five study areas.

	DMQC		RTQC – original		RTQC – first filter		RTQC – second filter	
	Meas.	%	Meas.	%	Meas.	%	Meas.	%
Alboran Sea	49 401	54.96 %	40 481	45.04 %	40 481	45.04 %	40 481	45.04 %
Antarctica	1 117 571	92.14 %	95 346	7.86 %	93 647	7.72 %	92 204	7.60 %
Bermuda Triangle	2 060 348	70.49 %	862 455	29.51 %	468 483	16.03 %	243 752	8.34 %
Hawaii	3 252 097	70.81 %	1 340 462	29.19 %	1 308 773	28.50 %	1 259 247	27.42 %
Indonesia	5 260 566	86.86 %	795 900	13.14 %	780 727	12.89 %	771 874	12.74 %

Table 3. Profilers and profiles present in the Mexican EEZ.

DAC	Core		Biogeochemical		Profiles
	Active	Inactive	Active	Inactive	
AO: AOML	51	114	0	3	32 998
IF: CORIOLIS	6	3	0	1	1098
ME: MEDS	1	1	0	0	201
Total	58	118	0	4	34 297

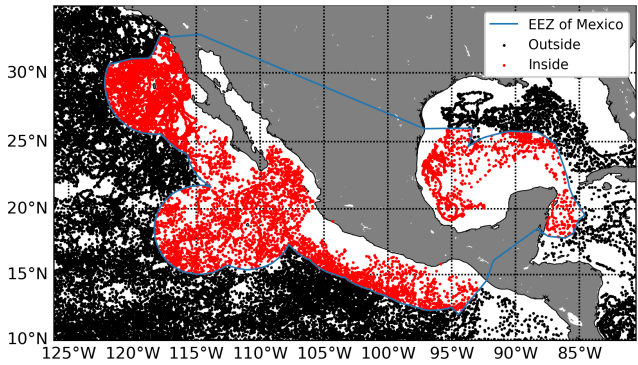


Figure 7. Filtered geographic locations within the EEZ of Mexico. The irregular polygon that delimits the EEZ of Mexico and the profiles measured inside and outside of it are shown.

gon given in the month of November 2019. We can see from the table that there is a shortage of biogeochemical profilers within the polygon. These four biogeochemical HAPs are capable of measuring oxygen in addition to temperature and salinity, but none of their oxygen data satisfactorily finish the quality control process, so they are not available. So we can conclude that within the Mexican EEZ there are no good-quality biogeochemical data from Argo HAPs.

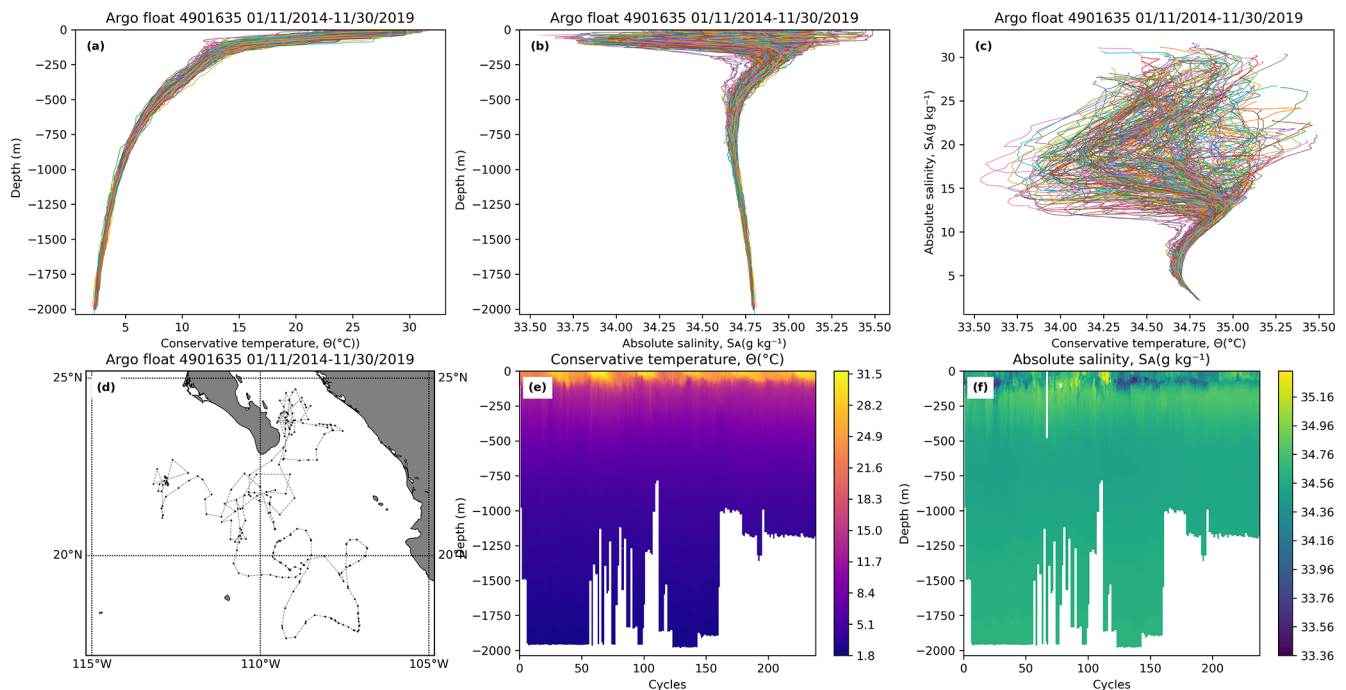


Figure 8. Diagrams produced by the web application. (a) Profile of conservative temperature. (b) Profile of absolute salinity. (c) T – S diagram. (d) HAP trajectory. (e) Profile of conservative temperature with respect to time. (f) Profile of absolute salinity with respect to time.

For each of these profilers their profiles of temperature (Fig. 8a) and salinity (Fig. 8b), the temperature–salinity (T – S) diagram (Fig. 8c), the estimation of the profiler trajectory (Fig. 8d), and the profiles of temperature (Fig. 8e) and salinity (Fig. 8f) with respect to time were generated. These diagrams are basic for analysis in scientific ocean research. The profiler 4901635 is shown as an illustrative example in Fig. 8.

The satellite map of the web application is interactive. It shows the active and inactive HAPs; filters the data; shows statistics, trajectories and diagrams (Fig. 9a); and has other tools to facilitate the visualization and management of the data, such as displaying statistics of a given study area within the main polygon (Fig. 9b and c).

Finally, the filtering of RTQC data that have patterns similar to DMQC data is offered in the web application, which allows the data to be filtered in a study area within the EEZ of Mexico. It is not necessary to have programming knowledge. Access to the web application is through the *cluster_qc* library repository.

4 Discussion

Despite the existence of reports on salinity drifts such as the one announced by Argo Data Management (2020) on 25 September 2018, the quality control processes in real time are not yet robust enough to identify them, since these processes are automatic and mainly look for impossible data, for example, latitudes greater than 90° or temperatures and

salinities outside the global and regional ranges. Therefore, the quality established by the flags does not take these drifts into account. A possible solution is for the interested user to apply the DMQC on their own (Wong et al., 2021). This process can be long and tedious, and for this reason and under Argo's recommendation to use only DMQC for scientific research, a large number of users who use the data from Argo HAPs prefer to directly discard the RTQC data and only use the DMQC data.

The data in DMQC are consistent with other international databases such as WOA18 within the study area delimited by the irregular polygon, which validates this process. However, too much data have to be discarded due to the drifts present in RTQC. The filtering proposed in this work is based on using the patterns followed by the DMQC data to filter the RTQC data, especially useful for areas where there are few profiles. This process is carried out by zone and by month, and in this way it does not matter if the study area is close to the arctic or the tropics. The filtering of the RTQC data are carried out based on the characteristics of the area reflected in the DMQC data. In addition, when separating the data by month, their seasonal changes are taken into account. This means that the resulting RTQC data will have a high probability of being accepted when the DMQC is applied to them.

The time it takes for a modern computer to do cluster analysis is relatively short compared to the 12 months it can take to perform the DMQC, and this will help users interested in recent data from HAPs to have greater reliability when using RTQC data. Two filters are proposed: the first is the result

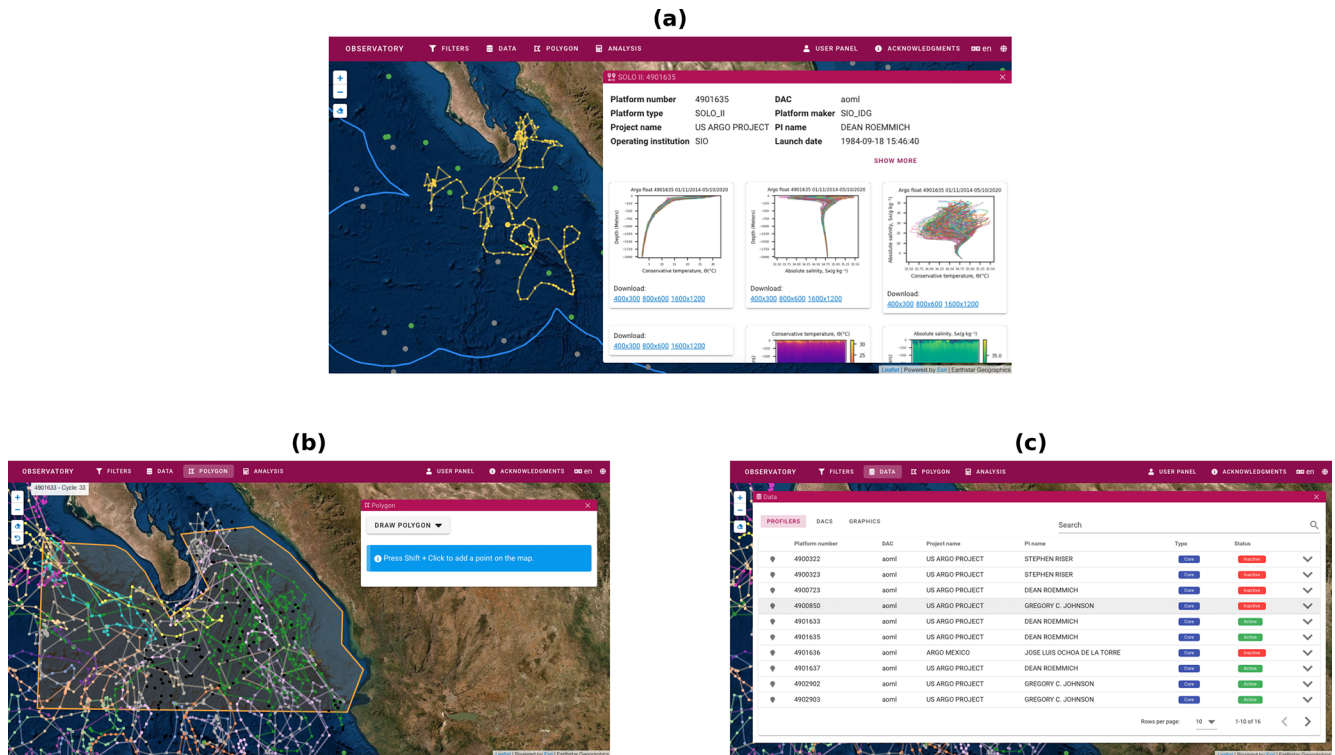


Figure 9. Web application. (a) Data, charts and trajectory of a HAP. (b) HAP trajectories filtered by a drawn polygon. (c) Profiler data within the polygon.

of using cluster analysis on the data and the second discards the HAPs that have presented salinity drifts in the result of the first filter. Therefore the second filter is more reliable but contains a smaller amount of data. As seen in the TPCM example, the user is free to use either one or a combination of both. For example, as seen in Fig. 6a and b, where around 80 % and 30 % of the total discarded data are admitted, the months of July to September continued with salinity drifts after applying the first filter, to take advantage of more data the user can use the data from the months of July to September of the second filter and the rest of the months use the data from the first. However, as shown by the five study areas used as an extra example, the percentage of data admitted by the filters depends on the study area and its characteristics. It is the responsibility of the users to make the decision based on their knowledge of the study area, which filter to use, if the study area should be resized or if the default depth value should be changed.

There are platforms to access data from Argo HAPs, such as Argo Data Management (2020), Coriolis (2020) and Euro-Argo (2020) in addition to other options such as FTP or snapshots. The current platforms already provide graphics and data from the profilers, as well as filters to display or download the data. However, the geographical filter they use is by maximum and minimum coordinates, so it is only possible to

filter by polygons in rectangle or square shape without rotation.

Another platform called OceanOPS (2021) (Joint Centre for Oceanography and Marine Meteorology in situ Observations Programmes Support) does perform statistical analyses on the data; nevertheless this one performs them globally, and it is not possible to choose a smaller area, for example, only the EEZ of Mexico or the tropical Pacific off central Mexico and surrounding areas, to obtain statistical information on it. It is worth mentioning that said platform has a large number of statistics for each variable registered within the source files; however, being able to generate graphs and tables in real time using an irregular polygon defined by the user (as shown in this work with the PIP algorithm) would be a great tool for studying these data.

The web application described in this document tries to cover some of the aforementioned problems and include some of their characteristics, in addition to proposing unpublished options such as filtering by irregular polygons, statistics adaptable to filters, generation of graphs according to user needs and RTQC data filtering. However, the web application is in its initial phase, and there are still many tools and databases that can be integrated to offer an even more complete experience.

5 Conclusions

This work gives two filtering methods to discard only the RTQC data that present salinity drifts and with it to take advantage of the largest amount of data within a given polygon. In the TPCM, from the total RTQC data it was possible to recover around 80 % in the case of the first filter and 30 % in the case of the second, which are usually discarded due to problems such as salinity drifts. This allows users to use any of the filters or a combination of both to have a greater amount of data within the study area of their interest in a matter of minutes, rather than waiting for the DMQC that takes up to 12 months to be completed.

This work provides useful tools to increase productivity in scientific investigations that use data from the water column. The PIP algorithm turns out to be an efficient method to directly filter the data from any georeferenced database using geographic locations, while the algorithms proposed for filtering RTQC data allow the separation of the data not yet adjusted by the DMQC into data with salinity drifts and data that show patterns similar to those of the DMQC data, in order to increase the amount of data in study areas with scarce data from HAPs. Finally, the web app demonstrates one of the applications in which these proposals can be used.

Code availability. *cluster_gc* was developed in Python 3.7 and is licensed under a Creative Commons Attribution 4.0 International License. The source code is available at <https://doi.org/10.5281/zenodo.4595802> (Romero et al., 2021a). The latest package version is v1.0.2.

Data availability. These data were collected and made freely available by the International Argo Program and the national programs that contribute to it: <https://argo.ucsd.edu> (last access: January 2020, Argo, 2020a) and <https://www.ocean-ops.org> (last access: April 2021, OceanOPS, 2021). The Argo Program is part of the Global Ocean Observing System. The data were downloaded from the Coriolis GDAC FTP server in 2019, and the snapshot from December 2020 (Argo, 2020b) was also used. The data used from the NCEI World Ocean Database 2018 are the monthly data of temperature (<https://www.ncei.noaa.gov/products/world-ocean-database>, last access: November 2019, Locarnini et al., 2018) and salinity (<https://www.ncei.noaa.gov/products/world-ocean-database>, last access: November 2019, Zweng et al., 2018) of the statistical mean of each quarter of a degree ($1/4^\circ$) from 2005 to 2017. Maps throughout this work were created using ArcGIS® software by Esri. ArcGIS® and ArcMap™ are the intellectual property of Esri and are used herein under license. Copyright © Esri. All rights reserved. For more information about Esri® software, please visit <https://www.esri.com/en-us/home> (last access: November 2019).

Author contributions. ER developed the methodology and software described in this work and also performed the data analysis. LTF

supervised this work. LTF, IC and MC contributed to the conceptualization and design of the study, the interpretation of the data, and the preparation of the article.

Competing interests. The authors declare that they have no conflict of interest.

Disclaimer. Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Acknowledgements. We are grateful to CONACYT for granting scholarship no. 731303 to Emmanuel Romero. We appreciate that these data were collected and made freely available by the International Argo Program and the national programs that contribute to it (<https://argo.ucsd.edu> (last access: January 2020) <https://www.ocean-ops.org>, last access: April 2021). The Argo Program is part of the Global Ocean Observing System. We also thank the Instituto Tecnológico de La Paz (ITLP) and the Centro Interdisciplinario de Ciencias Marinas (CICIMAR) for their institutional support. We also acknowledge the critical comments from the reviewers.

Financial support. This research has been supported by CONACYT (scholarship no. 731303).

Review statement. This paper was edited by Oliver Zielinski and reviewed by two anonymous referees.

References

- Argo: Argo [data set], available at: <https://argo.ucsd.edu/>, last access: January 2020a.
- Argo: Argo float data and metadata from Global Data Assembly Centre (Argo GDAC) – Snapshot of Argo GDAC of December 10st 2020, [data set], <https://doi.org/10.17882/42182#79118>, 2020b.
- Argo Data Management Team: Argo user's manual V3.3, Report, <https://doi.org/10.13155/29825>, 2019.
- Argo Data Management: Argo Data Management, available at: <http://www.argodatamgt.org/>, last access: 2020.
- Coriolis: Coriolis: In situ data for operational oceanography, available at: <http://www.coriolis.eu.org/>, last access: 2020.
- Euro-Argo: Argo Fleet Monitoring – Euro-Argo, available at: <https://fleetmonitoring.euro-argo.eu/>, last access: 2020.
- Everitt, B., Landau, S., Leese, M., and Stahl, D.: Cluster Analysis, Wiley Series in Probability and Statistics, Wiley, 346 pp., 2011.
- Fiedler, P. and Talley, L.: Hydrography of the Eastern Tropical Pacific: a review, Prog. Ocean., 69, 143–180, <https://doi.org/10.1016/j.pocean.2006.03.008>, 2006.
- Godínez, V. M., Beier, E., Lavín, M. F., and Kurczyn, J. A.: Circulation at the entrance of the Gulf of California from satellite al-

- timer and hydrographic observations, *J. Geophys. Res.-Oceans*, 115, C04007, <https://doi.org/10.1029/2009JC005705>, 2010.
- Hartigan, J. A. and Wong, M. A.: Algorithm AS 136: A K-Means Clustering Algorithm, *J. R. Stat. Soc.*, 28, 100–108, <https://doi.org/10.2307/2346830>, 1979.
- Foley, J. D., van Dam, A., Feiner, S. K., and Hughes, J. F.: *Computer Graphics: Principles and Practice*, The Systems Programming Series, Addison-Wesley, 2 Edn., 1175 pp., 1990.
- Kessler, W. S.: The circulation of the eastern tropical Pacific: A review, *Prog. Ocean.*, 69, 181–217, <https://doi.org/10.1016/j.pocean.2006.03.009>, 2006.
- Lavín, M., Beier, E., Gomez-Valdes, J., Godínez, V., and García, J.: On the summer poleward coastal current off SW México, *Geophys. Res. Lett.*, 33, L02601, <https://doi.org/10.1029/2005GL024686>, 2006.
- Lavín, M. F. and Marinone, S. G.: *An Overview of the Physical Oceanography of the Gulf of California*, Springer Netherlands, 173–204, https://doi.org/10.1007/978-94-010-0074-1_11, 2003.
- Lavín, M. F., Castro, R., Beier, E., Godínez, V. M., Amador, A., and Guest, P.: SST, thermohaline structure, and circulation in the southern Gulf of California in June 2004 during the North American Monsoon Experiment, *J. Geophys. Res.-Oceans*, 114, C02025, <https://doi.org/10.1029/2008JC004896>, 2009.
- Locarnini, R., Mishonov, A., Baranova, O., Boyer, T., Zweng, M., Garcia, H., Reagan, J., Seidov, D., Weathers, K., Paver, C., Smolyar, I., and Locarnini, R.: *World Ocean Atlas 2018, Volume 1: Temperature*, edited by: Mishonov, A., NOAA Atlas NESDIS [data set], available at: <https://www.ncei.noaa.gov/products/world-ocean-database> (last access: November 2019), 1, 52 pp., 2018.
- OceanOPS: *OceanOPS [data set]*, available at: <https://www.ocean-ops.org>, last access: April 2021.
- Pantoja, D., Marinone, S., Pares-Sierra, A., and Gomez-Valdivia, F.: Numerical modeling of seasonal and mesoscale hydrography and circulation in the Mexican Central Pacific, *Cienc. Mar.*, 38, 363–379, <https://doi.org/10.7773/cm.v38i2.2007>, 2012.
- Portela, E., Beier, E., Barton, E., Castro Valdez, R., Godínez, V., Palacios-Hernández, E., Fiedler, P., Sánchez-Velasco, L., and Trasviña-Castro, A.: Water Masses and Circulation in the Tropical Pacific off Central Mexico and Surrounding Areas, *J. Phys. Ocean.*, 46, 3069–3081, <https://doi.org/10.1175/JPO-D-16-0068.1>, 2016.
- Romero, E., Tenorio-Fernandez, L., Castro, I., and Castro, M.: *romeroqc/cluster_qc: Filtering Methods based on cluster analysis for Argo Data*, Zenodo [code], <https://doi.org/10.5281/zenodo.4595802>, 2021a.
- Romero, E., Tenorio-Fernandez, L., Castro, I., and Castro, M.: Argo data filtering results to avoid salinity drifts, <https://doi.org/10.6084/m9.figshare.14999613.v1>, 2021b.
- Stramma, L., Johnson, G. C., Sprintall, J., and Mohrholz, V.: Expanding Oxygen-Minimum Zones in the Tropical Oceans, *Science*, 320, 655–658, <https://doi.org/10.1126/science.1153847>, 2008.
- Wong, A., Keeley, R., and Carval, T.: *Argo Quality Control Manual for CTD and Trajectory Data*, Report, USA, France, <https://doi.org/10.13155/33951>, 2021.
- Zamudio, L., Leonardi, A., Meyers, S., and O'Brien, J.: ENSO and Eddies on the Southwest coast of Mexico, *Geophys. Res. Lett.*, 28, 2000GL011814, <https://doi.org/10.1029/2000GL011814>, 2001.
- Zamudio, L., Hurlburt, H., Metzger, E., and Tilburg, C.: Tropical Wave-Induced Oceanic Eddies at Cabo Corrientes and the Maria Islands, Mexico, *J. Geophys. Res.*, 112, 18, <https://doi.org/10.1029/2006JC004018>, 2007.
- Zweng, M., Reagan, J., Seidov, D., Boyer, T., Locarnini, R., Garcia, H., Mishonov, A., Baranova, O., Paver, C., and Smolyar, I.: *World Ocean Atlas 2018 [data set], Volume 2: Salinity*, available at: <https://www.ncei.noaa.gov/products/world-ocean-database> (last access: November, 2019), edited by: Mishonov, A., NOAA Atlas NESDIS, 50 pp., 2018.