



1 **Data mining-based machine learning methods for improving**
2 **hydrological data: a case study of salinity field in the Western**
3 **Arctic Ocean**

4

5 **Shuhao Tao^{1,2}, Ling Du^{1,2}, Jiahao Li^{1,2}**

6 *¹Frontier Science Center for Deep Ocean Multispheres and Earth System (FDOMES)*
7 *and Physical Oceanography Laboratory, Ocean University of China, Qingdao, China,*

8 *²College of Oceanic and Atmospheric Sciences, Ocean University of China, Qingdao,*
9 *China*

10

11 **Correspondence to:** Ling Du (duling@ouc.edu.cn)



12 **Abstract.** In the Western Arctic Ocean lies the largest freshwater reservoir in the Arctic
13 Ocean, the Beaufort Gyre. Long-term changes in freshwater reservoirs are critical for
14 understanding the Arctic Ocean, and data from various sources, particularly measured
15 or reanalyzed data, must be used to the greatest extent possible. Over the past two
16 decades, a large number of intensive field observations and ship surveys have been
17 conducted in the western Arctic Ocean to obtain a large amount of CTD data. Multiple
18 machine learning methods were evaluated and merged to reconstruct annual salinity
19 product in the western Arctic Ocean over the period 2003-2022. Data mining-based
20 machine learning methods make use of variables determined by physical processes,
21 such as sea level pressure, sea ice concentration, and drift. Our objective is to effectively
22 manage the mean root mean square error (RMSE) of sea surface salinity, which exhibits
23 greater susceptibility to atmospheric, sea ice, and oceanic changes. Considering the
24 higher susceptibility of sea surface salinity to atmospheric, sea ice, and oceanic changes,
25 which leads to greater variability, we ensured that the average root mean square error
26 of CTD and EN4 sea surface salinity field during the machine learning training process
27 was constrained within 0.25psu. The machine learning process reveals that the
28 uncertainty in predicting sea surface salinity, as constrained by CTD data, is 0.24%,
29 whereas when constrained by EN4 data it reduces to 0.02%. During data merging and
30 post-calibrating, the weight coefficients are constrained by imposing limitations on the
31 uncertainty value. Compared with commonly used EN4 and ORAS5 salinity in the
32 Arctic Ocean, our salinity product provide more accurate descriptions of freshwater
33 content in the Beaufort Gyre and depth variations at its halocline base. The application
34 potential of this multi-machine learning results approach for evaluating and integrating
35 extends beyond the salinity field, encompassing hydrometeorology, sea ice thickness,
36 polar biogeochemistry, and other related fields. The datasets are available at
37 <https://zenodo.org/records/10990138> (Tao and Du, 2024).

38

39 1. Introduction

40 Unlike the low- and mid-latitude oceans, the Arctic Ocean is characterized by its
41 extensive sea ice coverage and near-freezing sea surface water. Variations in salinity in
42 the Western Arctic Ocean have profound implications for stratification strength, ocean
43 circulation patterns, and biogeochemical cycles (Carmack et al., 2016; Cornish et al.,
44 2020). Freshwater reservoirs and their evolution, which are closely related to the change
45 of seawater salinity, have become the focus of research in the Arctic Ocean. Therefore,
46 obtaining accurate salinity data holds great significance for our understanding of this
47 unique marine environment. The mean density structure and wind-driven surface



48 circulation in the Arctic Ocean are predominantly influenced by two key factors: The
49 anti-cyclonic Beaufort Gyre located in the Canadian Basin and the Transpolar Drift
50 ([Hall et al., 2022](#)). Furthermore, within Western Arctic Oceans, significant amounts of
51 freshwater accumulate within the Beaufort Gyre. The release of this freshwater exerts
52 a substantial impact on local climate dynamics as well as global climate change at large
53 scales ([Carmack et al., 2008](#); [Giles et al., 2012](#); [Proshutinsky et al., 2009, 2019](#)). Our
54 research specifically focuses on a case study of investigating salinity product improved
55 by multi-machine learning results evaluating and integrating within Western Arctic
56 Oceans.

57 The presence of sea ice severely limits the availability of salinity data in the Arctic
58 Ocean, posing significant challenges to meeting the demands of current research.
59 Shipborne observations of CTD and ITP data are sporadic, posing challenges in
60 obtaining reliable salinity measurements. The accuracy of both model and reanalysis
61 data is frequently subpar. [Behrendt et al. \(2018\)](#) collected a large amount of measured
62 data to form a Unified Database for Arctic and Subarctic Hydrography for the period
63 1980-2015, however, hydrological data for recent years are lacking.. In recent years,
64 however, highly developed measurement techniques were especially designed for
65 operation in the Arctic environment. Furthermore, an increasing number of research
66 activities and international collaboration - such as Beaufort Gyre Exploration Project
67 (BGEP) has generated a large number of hydrographic data in the Western Arctic ocean
68 and the subarctic seas (e.g., [Rabe et al., 2014](#)).

69 The advancement of stochastic computer science and technology in recent years has led
70 to an increasing utilization of machine learning methods across various domains. The
71 utilization of data mining-based machine learning techniques for data generation is
72 explored in this paper, with a focus on the salinity observed in the Western Arctic Ocean.
73 Machine learning techniques have already demonstrated their efficacy in data
74 generation tasks. For instance, [Wang et al. \(2023\)](#) employed a machine-learning-based
75 regression method to reconstruct long-term (2003-2020) sea surface pCO₂ in the South
76 China Sea, while [Chen et al. \(2024\)](#) utilized the Random Forest Algorithm to generate
77 datasets of stable isotopes of precipitation in the Eurasian continent. The utilization of
78 machine learning offers distinct advantages during data reconstruction processes
79 including high automation, exceptional accuracy, robust scalability, and expedited
80 processing compared to assimilation approaches. Consequently, this paper employs
81 several machine learning methods to produce dependable salinity data in the western
82 Arctic Ocean.

83 We performed machine learning training on sea level pressure, sea ice concentration,

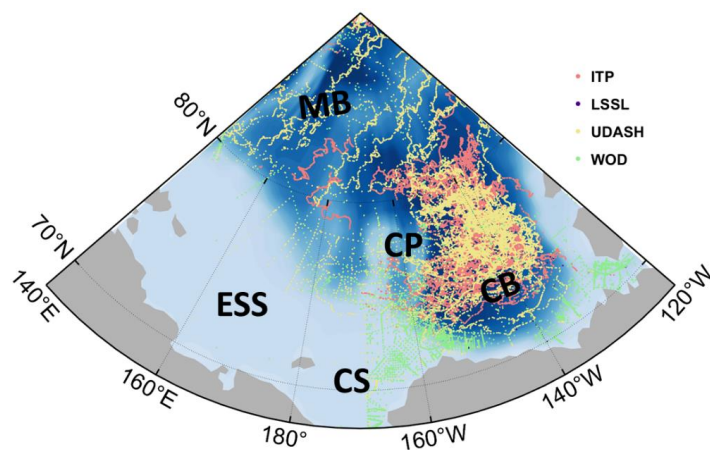


84 sea ice motion, as well as a large number of quality-controlled CTD data and EN4 data
85 using various machine learning methods. The datasets were merged to generate a
86 salinity product with a resolution of $0.5 \times 0.25^\circ$ above 1000m for the period spanning
87 from 2003 to 2022, encompassing a total of 48 vertical layers. The machine learning
88 performance was assessed not only through RMSE, but also by evaluating the
89 uncertainty resulting from data merging and post-calibrating processes. The ORAS5
90 and EN4 datasets were employed to investigate the Beaufort Gyre and Arctic Ocean
91 (Hall et al., 2022). The accuracy and reliability of our salinity product were
92 demonstrated by comparing it with EN4 and ORAS5 data, as well as measured
93 freshwater content and halocline base depth in the Beaufort Gyre region.

94 **2. Data and methodology**

95 **2.1 Study area**

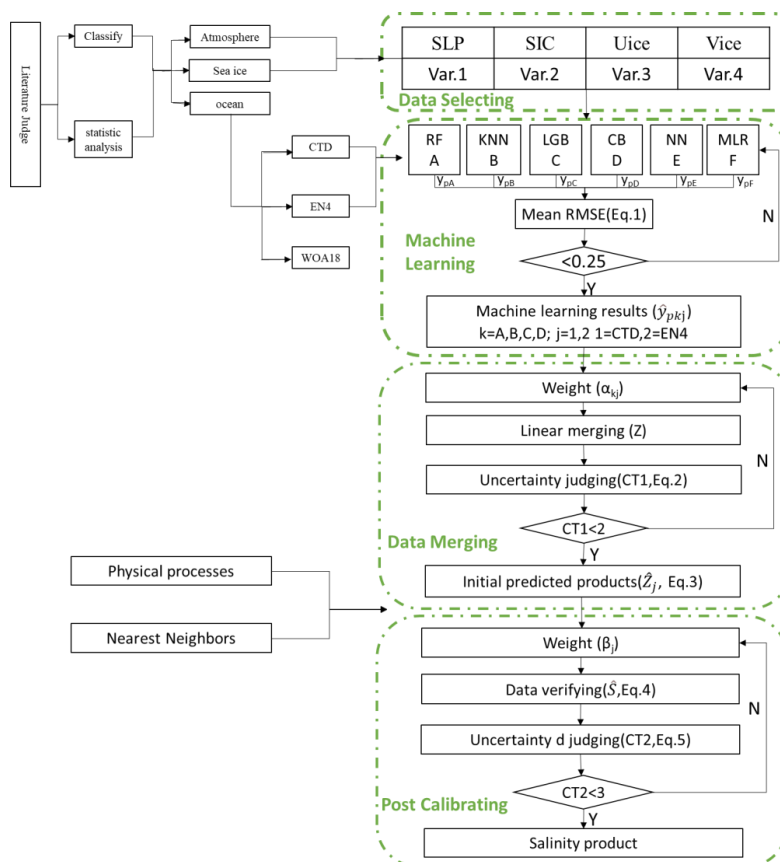
96 The Western Arctic Ocean (140°E - 120°W , 68°N - 90°N) spans a vast territory with the
97 Beaufort Gyre, the largest fresh water reservoir in the Arctic Ocean (Fig. 1). In the
98 Western Arctic Ocean, sea ice covers the area in winter, while in summer, a large area
99 of sea ice at low latitudes melts. However, sea ice still exists in the multi-year ice zone
100 in the northeast of Canada Basin. The Western Arctic Ocean is mainly influenced by
101 the anticyclonic Beaufort High. In the western part of the Arctic Ocean, there is the
102 main circulation system of the Arctic Ocean, the Beaufort Gyre, which accumulates a
103 large amount of fresh water. The Strength of the Beaufort Gyre has been continuously
104 increasing, reaching a stable state after 2007, with changes in freshwater content
105 consistent with the strength of the gyre (Regan et al., 2019). The range of the Beaufort
106 Gyre expanded westward from 2003 to 2013, and contracted eastward back to the
107 Canadian Basin after 2014 (Lin et al., 2023). Freshwater accumulation, storage, and
108 release from the BG exert far-reaching impacts on both regional and global climate
109 systems. Therefore, accurate salinity data is very important for our study of Beaufort
110 Gyre.



111

112 **Figure1. Topography of the Western Arctic Ocean. The map also includes the**
113 **Canada Basin (CB), Chukchi sea (CS), the Chukchi Plateau (CP), East Siberian**
114 **Sea (ESS) and Makarov Basin (MB).**

115 Our goal is to generate a set of salinity product that can be used to analyze the physical
116 ocean environment changes in the Arctic Ocean in recent years. The procedure of
117 improving salinity product is mainly divided into four major parts, which are data
118 selecting, machine learning training, data merging, and post-calibrating.



119

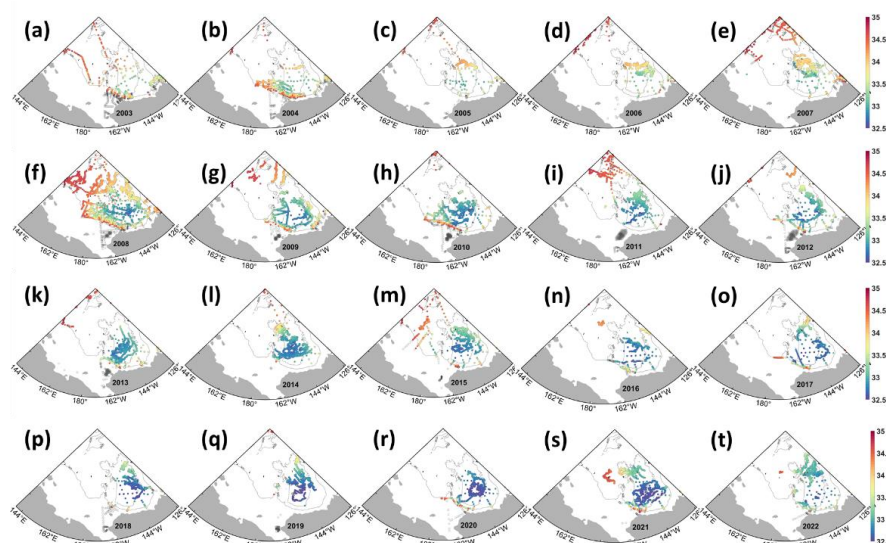
120 **Figure.2 Procedure for improving the salinity field in the Western Arctic Ocean**
 121 **through a data mining-based machine learning method**

122 **2.2 Data Selecting**

123 We have collected a large amount of CTD salinity data. The World Ocean Database
 124 (WOD) is world's largest collection of uniformly formatted, quality controlled, publicly
 125 available ocean profile data ([https://www.ncei.noaa.gov/access/world-ocean-](https://www.ncei.noaa.gov/access/world-ocean-database/bin/getwodyearlydata.pl?Go=TimeSorted)
 126 [database/bin/getwodyearlydata.pl?Go=TimeSorted](https://www.ncei.noaa.gov/access/world-ocean-database/bin/getwodyearlydata.pl?Go=TimeSorted), last access: 8 December 2023). We
 127 selected the WOD18 salinity profiles and retained the data with flags 0 and 1 based on
 128 the quality control provided by the data itself. Unified Database for Arctic and Subarctic
 129 Hydrography (UDASH) is a unified and high-quality temperature and salinity data set
 130 for the Arctic Ocean and the subpolar seas north of 65° N for the period 1980-2015
 131 (<https://essd.copernicus.org/articles/10/1119/2018/>, last access: 8 December 2023). Sea
 132 ice presents a significant impediment to sustained observation of the Arctic Ocean.



133 Researchers designed and field tested an automated, easily-deployed Ice-Tethered
134 Profiler (ITP) for Arctic study. Building on the ongoing success of ice drifters that
135 support multiple discrete subsurface sensors on tethers and the WHOI-developed
136 Moored Profiler instrument capable of moving along a tether to sample at better than
137 1-m vertical resolution (<https://www2.whoi.edu/site/itp/data/>, last access: 8 December
138 2023). Shipboard hydrographic data and water sampling measured on board the CCGS
139 Louis S. St-Laurent (LSSL) are carried out at about 30 standard sites on each cruise
140 (<https://www2.whoi.edu/site/beaufortgyre/data/ctd-and-geochemistry/>, Last access: 8
141 December 2023), the CTD data of LSSL collected during the 2004 expedition was not
142 utilized.



143

144 **Figure 3. Annual sea surface salinity fields from 2003 to 2022 in the Western Arctic**
145 **Ocean.**

146 The data collected include a variety of issues such as missing values, outliers, and
147 duplicates as well as gaps in dates and missing or incorrect latitude and longitude
148 information. Therefore, the collected raw data underwent pre-processing and data
149 cleaning. Missing data were interpolated, entries that could not be completed were
150 removed, and duplicate data were eliminated. This article interpolates all data onto the
151 WOD vertical grid in depth. The most CTD data was collected in late summer and early
152 autumn (August to October), while the least CTD data was collected in June. The
153 measured data is mainly concentrated in the Canadian Basin, with very few measured



154 data in the East Siberian Sea (Fig. 2). After 2003, ITP and LSSL supplemented a large
155 amount of CTD data in situ, so we hope to generate gridded data from 2003 to 2022.

156 In addition to a large amount of observed CTD data, considering the temporal and
157 spatial discontinuity of the observed data, we have introduced EN4
158 (<https://www.metoffice.gov.uk/hadobs/en4/>, Last access: 8 December 2023) reanalysis
159 data. Furthermore, taking into account the influence of the atmosphere and sea ice on
160 the ocean, we have also incorporated SLP data from ERA5
161 ([https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels-
162 monthly-means?tab=form](https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels-monthly-means?tab=form), Last access: 8 December 2023) and sea ice concentration
163 and sea ice drift field data from NSIDC (<https://nsidc.org/home>, Last access: 8
164 December 2023). We use monthly salinity data provided by the European Centre for
165 Medium-Range Weather Forecasts (ECMWF) through the Ocean Reanalysis System's
166 version 5 (ORAS5), which uses the Nucleus for European Modeling of the Ocean
167 (NEMOv3.4) for its ocean model coupled with a sea ice model to assess the accuracy
168 of salinity product. In the data selecting section, we summarized previous literature and
169 selected the sea level pressure field, sea ice concentration, and sea ice drift field data of
170 the Western Arctic Ocean as training variables for machine learning.

171 **2.3 Machine learning**

172 In the second part of the machine learning training section, we selected six commonly
173 used machine learning methods, which are Random Forest (RF), K Nearest Neighbor
174 (KNN), LightGBM (LGB), CatBoost (CB), Neural Network (NN), and Multilinear
175 Regression (MLR). We determined the optimal value of different machine learning
176 algorithm using optuna hyper parameter methods (code from
177 <https://github.com/optuna/>, last access: 20 March 2024) and GridSearchCV (from
178 scikit-learning) for the training set. We trained EN4 and CTD data with six different
179 machine learning methods respectively.

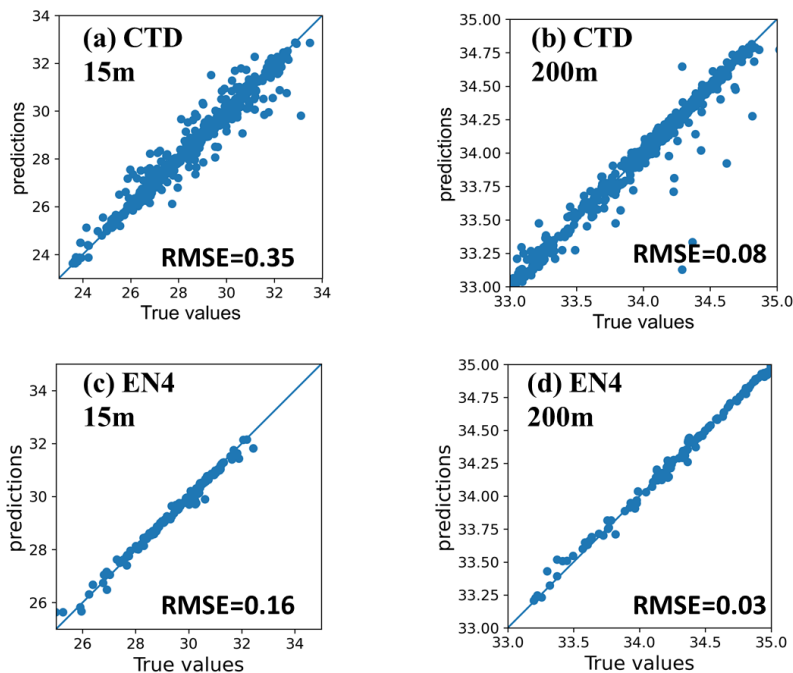
180 It is necessary to evaluate the accuracy of any model based on certain error metrics
181 before applying it to specific scenarios. Common model evaluation metrics include
182 MAE, RMSE. The mean squared error (MSE) is the standard deviation of the residuals
183 (prediction error), and the residuals are the distances between the fitted line and the data
184 points (i.e., the residuals show the degree of concentration of the reconstructed data
185 around the regression line). In regression analysis, RMSE is commonly used to verify
186 experimental results. To assess bias, the RMSE needs to combine the magnitude of the
187 model data and is calculated as follows:



$$188 \quad RMSE_{Kj} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_{iKj} - y_{piKj})^2}, \quad (Eq. 1:)$$

189 where n is the number of data points; K represents different machine learning
190 algorithms, and there are six types in total, which are RF, KNN, LGB, CB, NN, MLR;
191 $j=1$ represents CTD data, $j=2$ represents EN4 data; y is the training target data; y_p is
192 the prediction result after machine learning training.

193 Taking the results from 2008 of Random Forest results as an example (Fig.4), we found
194 that the salinity prediction at a depth of 200m is better than the prediction at the surface
195 (15m), and the prediction using EN4 data is better than using CTD data. However, what
196 is exciting is that even for the weakest prediction ability of CTD at the surface, the
197 RMSE is less than 0.35psu. Therefore, our evaluation of the model learning results will
198 mainly focus on the surface with larger prediction errors by RMSE.



199

200 **Figure 4. Comparisons between the predicted salinity and train target salinity**
201 **values for the Random Forest testing pool in 2008.**

202 In addition to RF, we also evaluated the prediction results of surface salinity for five



203 other machine learning methods using RMSE (Table1), which is calculated as follows:

204 **Table1. Evaluation of predicted surface salinity using different machine learning**
205 **methods**

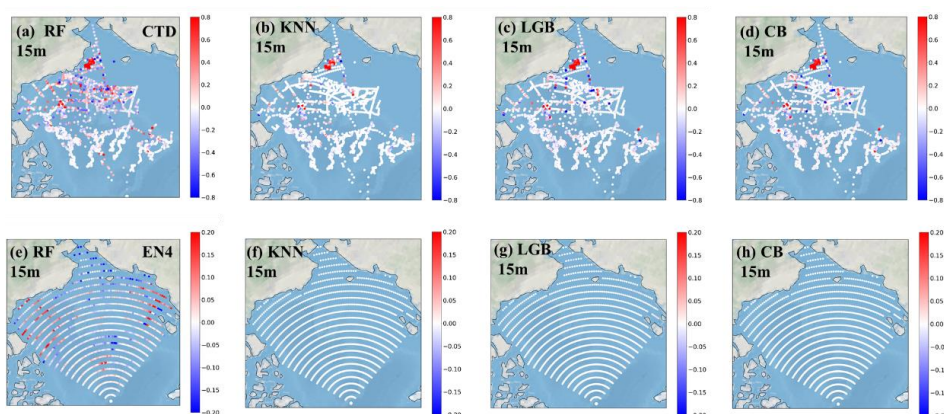
	Random Forest		K Nearest Neighbor		LightGBM		Catboost		Multilinear Regression		Neural Network	
	CTD	EN4	CTD	EN4	CTD	EN4	CTD	EN4	CTD	EN4	CTD	EN4
2003	0.45	0.07	0.49	0	0.43	0.00	0.43	0.00	1.07	0.90	1.01	0.52
2004	0.28	0.06	0.22	0	0.17	0.00	0.17	0.00	1.13	0.92	0.96	0.46
2005	0.08	0.09	0.09	0	0.06	0.00	0.08	0.00	0.57	0.97	0.34	0.55
2006	0.11	0.07	0.15	0	0.12	0.00	0.12	0.00	0.72	0.90	0.44	0.45
2007	0.19	0.06	0.21	0	0.18	0.00	0.19	0.00	1.14	1.10	0.79	0.41
2008	0.21	0.08	0.26	0	0.22	0.00	0.21	0.00	1.18	1.03	0.77	0.61
2009	0.12	0.07	0.16	0	0.13	0.00	0.12	0.00	0.82	0.99	0.52	0.63
2010	0.23	0.11	0.31	0	0.22	0.01	0.22	0.00	1.00	1.08	0.61	0.66
2011	0.17	0.10	0.22	0	0.17	0.00	0.16	0.00	1.00	0.92	0.54	0.57
2012	0.24	0.10	0.30	0	0.25	0.01	0.24	0.01	0.70	0.91	0.49	0.69
2013	0.20	0.08	0.28	0	0.20	0.00	0.20	0.00	0.70	0.86	0.45	0.54
2014	0.15	0.07	0.19	0	0.15	0.00	0.15	0.00	0.43	0.94	0.35	0.51
2015	0.18	0.07	0.21	0	0.17	0.00	0.17	0.00	0.61	0.87	0.48	0.48
2016	0.09	0.07	0.04	0	0.04	0.01	0.04	0.00	0.43	1.01	0.34	0.45
2017	0.21	0.09	0.04	0	0.06	0.00	0.04	0.00	0.68	0.91	0.57	0.55
2018	0.14	0.07	0.15	0	0.15	0.00	0.15	0.00	0.51	0.87	0.34	0.54
2019	0.34	0.06	0.28	0	0.25	0.00	0.19	0.00	1.00	0.89	0.78	0.56
2020	0.53	0.10	0.90	0	0.28	0.00	0.27	0.00	0.89	0.94	0.67	0.61
2021	0.38	0.07	0.45	0	0.34	0.00	0.13	0.00	0.88	0.82	0.76	0.53
2022	0.26	0.08	0.34	0	0.27	0.00	0.26	0.00	0.82	0.93	0.63	0.60

206

207 We selected four machine learning methods that prediction is closer to the training
208 target of sea surface salinity (with the mean RMSE less than 0.25), which are RF, KNN,
209 LGB, and CB. These four machine learning methods have better prediction results for
210 EN4 than for CTD. The errors generated during the prediction process mainly come
211 from the prediction of CTD salinity. The annual differences in predictive capabilities of
212 these four types of machine learning are very significant. The prediction results for RF
213 were the best in 2005 and 2016, and the worst in 2020, KNN had the best prediction
214 results for 2016 and 2017, and the worst prediction results for 2020. LGB had the best
215 forecast results for 2016 and 2017, and the worst forecast results for 2003. CB had the
216 best forecast results for 2016 and 2017, and the worst forecast results for 2003. In the
217 same year, some machine learning predictions are good while others are poor. For
218 example, in 2020, the predictions of RF and KNN were poor, but the predictions of
219 LGB and CB were good. This indicates that using multiple machine learning methods
220 can help improve the predictions of a certain method that performed poorly in a
221 particular year, eliminate biases in selecting machine learning methods for predictions,
222 and make the predictions more reliable.



223 RMSE is the spatial average result (Table 1), so only considering the numerical value
224 of RMSE will ignore the predictive ability of machine learning methods on different
225 regions in space. After training, we selected four machine learning methods with the
226 mean RMSE less than 0.25, which are RF, KNN, LGB, and CB. We take the example
227 of the prediction error of surface salinity in 2008 (predicted value minus training target
228 value) to analyze the salinity prediction ability of machine learning methods in different
229 regions. Machine learning models has significant spatial differences in predicting
230 salinity of CTD. Specifically, there are larger prediction errors in the Chukchi Sea,
231 Chukchi Sea Shelf, southern continental shelf slope of the Beaufort Gyre and center
232 Canada basin. The largest error occurred in the Chukchi Sea, which may be due to the
233 influence of Pacific water on the salinity of the upper layer of the Western Arctic Ocean.
234 The four machine learning methods for predicting surface salinity in EN4 are all very
235 good. KNN, LGB, and CB even have negligible prediction errors. RF shows a
236 significant spatial distribution in predicting surface salinity in EN4, with
237 overestimations in the southeast of the Canadian Basin and the western part of the East
238 Siberian Sea, with prediction errors less than 0.2psu. The predictions are
239 underestimated in the Chukchi Sea and the East Siberian Sea. The prediction errors of
240 different machine learning methods vary, so different weights need to be considered in
241 the data merge process.



242

243 **Figure 5. Error between the predicted salinity and real salinity values for the**
244 **training pool in 2008.**

245 **2.4 Data merging and post-calibrating**

246 The third part is the data merge part, where we linearly merging the training results
247 of the four better machine learning models. MAE is the average absolute difference



248 between the in situ data (true values) and the model output (predicted values). The sign
249 of these differences is ignored so that cancelations between positive and negative values
250 do not occur. RMSE and MAE have primarily been used to represent the uncertainties
251 in reconstructed datasets. In this article, we choose MAE as the criterion for assessing
252 uncertainty. We introduced weights and defined uncertainty, with uncertainty less than
253 2% as the indicator for selecting weights a_{kj} . The uncertainty (CT1) is calculated as

254 follows: $CT_{1kj} = \frac{1}{4} \sum_{k=1}^4 \frac{|\hat{y}_{kj} - Z_j|}{Y_j} \times 100\%$ (Eq. 2), Where k represents different

255 machine learning algorithms, and there are six types in total, which are RF, KNN, LGB,
256 CB; j=1 represents CTD data, j=2 represents EN4 data; $Z_j = \sum_{k=1}^4 a_{kj} \hat{y}_{kj}$ (Eq.3).

257 From this, we obtain the initial predicted products.

258 The salinity product is generated through the fourth post-calibrating, when there are
259 CTD measured data around the grid point, the salinity value of the point is formed by
260 merging the EN4 prediction results and the CTD prediction results according to weights;
261 otherwise, the salinity value of the point is taken as the EN4 prediction result. We
262 introduced weights and defined uncertainty, with uncertainty less than 3% as the
263 indicator for selecting weights β_{kj} . We need to check that salinity product $\hat{S} =$

264 $\sum_{j=1}^2 \beta_j \hat{Z}_j$ (Eq. 4) by uncertainty judging. The uncertainty (CT2) is calculated as

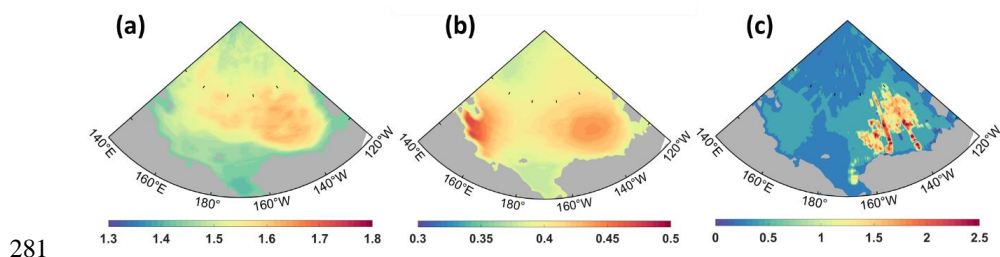
265 follows: $CT_{2j} = \frac{1}{2} \sum_{j=1}^2 \frac{|\hat{Z}_j - \hat{S}|}{\hat{S}} \times 100\%$ (Eq.5), Where j=1 represents CTD data, j=2

266 represents EN4 data. From this, we obtain the final salinity product in the Western
267 Arctic Ocean.

268 The uncertainty of the data in this article (represented by rMAE) includes three parts:
269 one part is the uncertainty generated during the machine learning process, with an
270 uncertainty of 0.24% for the surface salinity prediction generated by CTD and 0.02%
271 for the surface salinity prediction generated by EN4; the other parts include
272 uncertainties in data merging (Fig. 6a, 6b) and post calibrating (Fig. 6c). There are two
273 sets of initial predicted products for data merging of machine learning methods, EN4
274 and CTD. The uncertainty generated shows that the uncertainty constrained by CTD
275 data is larger in the central part of the Canadian Basin and the Chukchi Sea Shelf and
276 its adjacent waters, reaching 1.63% in the central part of the Canadian Basin. The
277 uncertainty constrained by EN4 data is larger in the central part of the Canadian Basin
278 and the East Siberian Sea, reaching 0.44% in the East Siberian Sea. The uncertainty
279 generated during the post-calibrating process is highest in the Canadian basin, with a



280 maximum value of 2.54%.

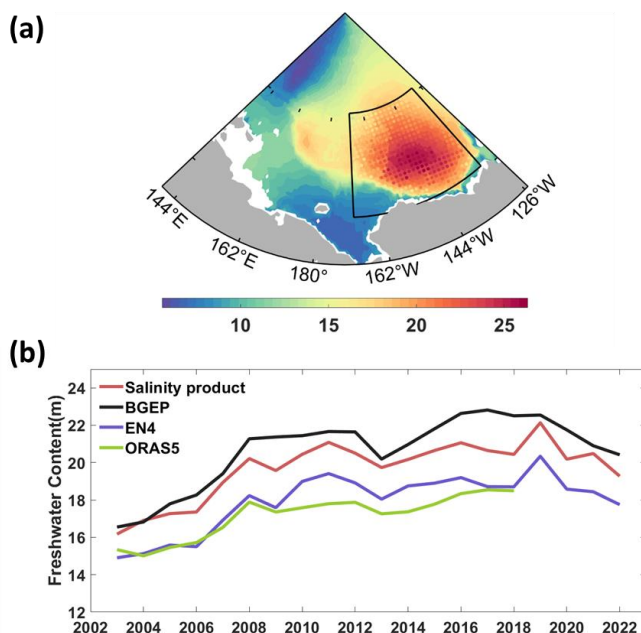


281
282 **Figure 6. Spatial pattern of sea surface salinity uncertainty (%) during the data**
283 **merging (a, CTD; b, EN4) and post-calibrating.**

284 3. Result and Discussion

285 We used the salinity product to calculate the freshwater content in the Beaufort Gyre
286 region (black box in Fig. 7a). In order to verify the superiority of the generated salinity
287 data in calculating the freshwater content, in addition to the freshwater content data
288 provided by BGEP for verification. On the other hand, the research of [Hall et al. \(2022\)](#)
289 showed that the salinity of ORAS5 and EN4 can be used to calculate the freshwater
290 content of the Arctic Ocean, and we also introduced the results of the freshwater content
291 calculation of ORAS5 (Fig. 7b). The FWC was computed relative to salinity 34.8 psu
292 following [Proshutinsky et al. \(2009\)](#):
$$FWC = \int_{z_{34.8}}^{z_{surface}} \left(\frac{34.8 - s(z)}{34.8} \right) dz \quad (Eq. 6)$$

293 The absolute errors of the freshwater content calculated by the generated salinity
294 product, the salinity data of EN4 and ORAS5 and the freshwater content provided by
295 BGEP are 4.89%, 13.21% and 16.40%, respectively. Using the generated salinity
296 product to calculate the freshwater content in the Beaufort Gyre region area can
297 improve the accuracy. We compared the spatial distribution of freshwater content
298 calculated from salinity product with freshwater content provided by BGEP. There are
299 areas on the Mendeleev Ridge with large freshwater content, which may be formed by
300 fresh water advection from the East Siberian Sea or by freshwater advection from the
301 Beaufort Gyre.



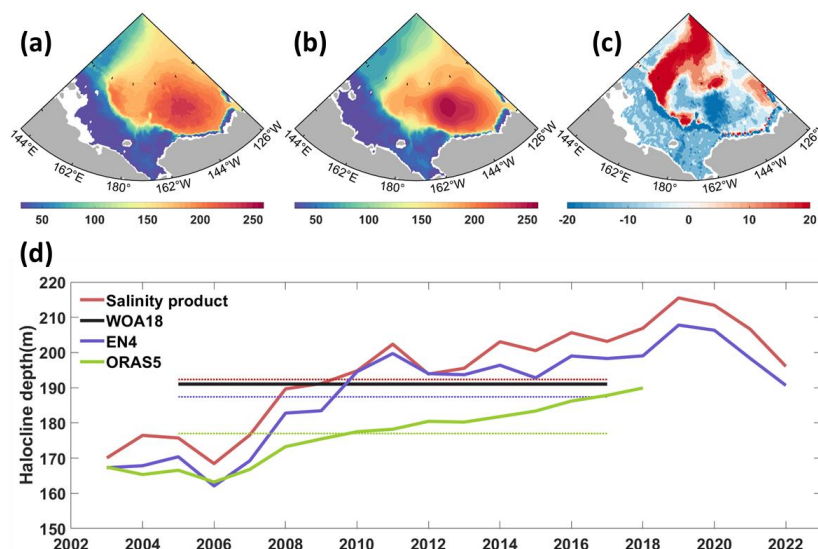
302

303 **Figure 7. Temporal and special variation of Freshwater Content (FWC, m).**
304 **(a)Shadow is Mean FWC from 2003 to 2022 derived from salinity product, color**
305 **dots represent FWC provided by BGEP. (b) Time series of FWC in Beaufort Gyre**
306 **region, Beaufort Gyre region is the black box in (a).**

307 The depth of halocline base plays an important role in studying the Beaufort Gyre
308 dynamics (e.g. Manucharyan et al.,2016). The depth of the halocline base is determined
309 by taking the 33.9 psu isosalinity line (Lin et al.,2023; Nyugen et al.,2012). All salinity
310 data used were interpolated vertically to 2m to calculate the depth of the halocline base.
311 The salinity product, EN4, ORAS5 and WOA18 calculated the halocline base depth in
312 Beaufort Gyre region of 192m,191m,187m and 176m, respectively (Fig. 8d). Salinity
313 product allow more accurate calculation of depth of halocline depth. Compared with
314 the results of ORAS5, the depth of halocline calculated by salinity product increased
315 significantly in the 2000s. Compared with EN4 results, the deepening trend in the 2010s
316 is more significant, but smaller than that of ORAS5. We compared the spatial
317 distribution characteristics of the bottom halocline and WOA18 obtained from salinity
318 product. The depth of halocline base is the deepest in the Canadian Basin, but the
319 salinity product results are shallower and more easterly than WOA18. The depth of the
320 halocline base calculated by salinity product is obviously 21m shallower in the
321 southwest of the Canadian Basin and 23m deeper in the north of the East Siberian Sea



322



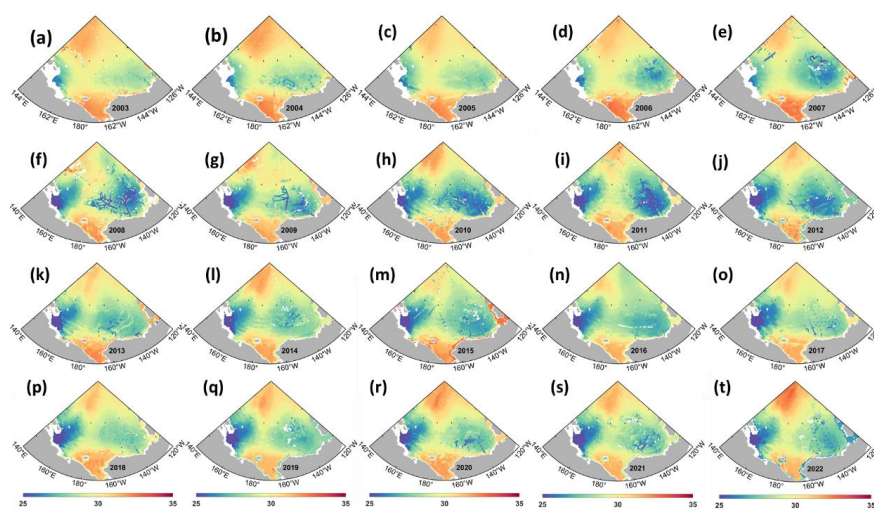
323

324 **Figure8. Temporal and special variation of Halocline depth (m). (a)Mean halocline**
325 **depth from 2005 to 2017 derived from salinity product (b) Mean halocline base**
326 **depth from 2005 to 2017 derived from salinity of WOA18. (c) Mean halocline**
327 **depth difference between salinity product and WOA18 from 2005 to 2017. (d)**
328 **Time series of halocline depth in Beaufort Gyre region.**

329 The results of salinity product indicate that the surface salinity is characterized by low
330 salinity in the central Canadian Basin and the East Siberian Sea, which indicates the
331 accumulation of fresh water there (Fig. 9). The continuous decrease in surface salinity
332 before 2011 and the continuous increase in surface salinity after 2011 indicate that
333 freshwater accumulated mainly at the surface before 2011 and decreased after 2011,
334 which support the recent major freshening event from 2012 to 2016 in North Atlantic
335 (Holliday et al.,2020). In the east-west direction, the surface low salt characteristics
336 westward expanded from 2003 to 2013, and eastward from 2014 to 2022, which
337 supports the conclusion that Beaufort Gyre expands westward (Regan et al.,2019;
338 Armitage et al., 2017) and shrinks eastward (Lin et al.,2023). In the north-south
339 direction, the surface low salt characteristics expanded northward in 2007, 2008, 2015
340 and 2016. The surface salinity of the East Siberian Sea decreased significantly in 2008
341 and has remained at reduced levels since then. According to the characteristics of
342 surface ocean circulation (Armitage et al., 2017), surface freshwater in the East Siberian



343 Sea may enter the Beaufort Gyre or flow out of the Arctic Ocean along the transpolar
344 drift. The characteristics of sea surface salinity can be seen that the Pacific water flows
345 partly to the northern Chukchi Sea, partly to the Canadian Basin and partly to the CAA
346 along the Alaskan coastal current, the reduced sea surface salinity of the Alaskan coastal
347 current indicates that less Pacific water is being transported along this path, indicating
348 a weakening of the Alaskan coastal current, whether this is influenced by the enhanced
349 Beaufort Gyre.



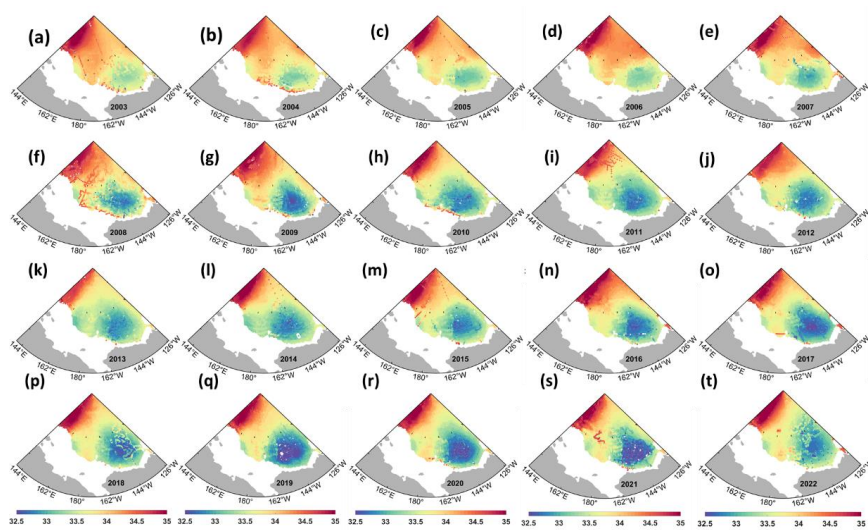
350

351 **Figure 9. Annual sea surface salinity fields in the Western Arctic ocean from 2003**
352 **to 2022. The color dots represent the measured CTD results, and the white dots**
353 **represent the measured sites that were deleted after quality control (see section**
354 **2.2).**

355 In order to observe the salinity distribution at the bottom of the halocline, which is about
356 200m deep in the western Arctic Ocean, we have analyzed the salinity distribution at
357 200m (Fig. 10). The results of salinity product indicate that salinity at 200m is
358 characterized by low salinity in the central Canadian Basin which indicates the
359 accumulation of fresh water in Canada Basin. Unlike the sea surface salinity, the salinity
360 at 200m has remained a slow downward trend after a rapid decline before 2008. This
361 suggests that fresh water in the Canadian Basin was relatively stable after a rapid
362 accumulation prior to 2008. Prior to 2008, freshwater in the western Arctic Ocean
363 pooled in large quantities at both the surface and the bottom of the halocline. After 2008,
364 the surface water decreased significantly while the bottom of the halocline water still
365 increased, indicating that the freshwater may be redistributed in the Arctic Ocean



366 through westward and northward expansion into the Marklov Basin (Bertosio et
367 al.,2022) or transported out of the Arctic Ocean (Zhang et al.,2021), or it may be pooled
368 deeper into the water column. From 2003 to 2013, the range of low salinity
369 characteristics of the halocline depth expanded, indicating that the area of freshwater
370 reservoir expanded and the area of Beaufort Gyre expanded. The salinity at 200m in
371 2022 increases significantly, indicating that there may be a freshwater migration
372 process in 2022.



373

374 **Figure 10. Reconstructed annual salinity fields at 200m in the Western Arctic**
375 **ocean from 2003 to 2022.**

376 4. Data availability

377 The salinity product ($0.5 \times 0.25^\circ$, 2003-2022) is available at
378 <https://zenodo.org/records/10990138> (Tao and Du, 2024).

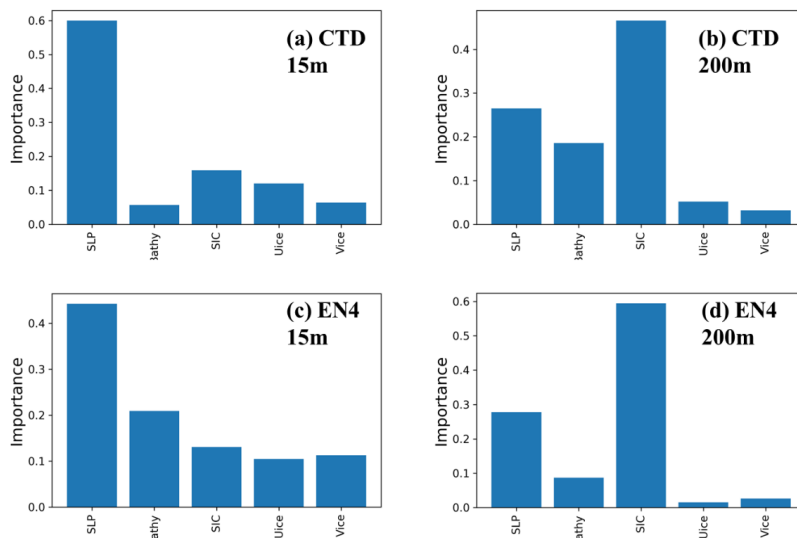
379 5. Summary

380 Based on data mining-based machine learning method, we have provided a salinity
381 product for the Western Arctic Ocean with a resolution of $0.5^\circ \times 0.25^\circ$ for the period
382 spanning from 2003 to 2022. This was achieved by establishing correlations between
383 bathymetry, sea ice dynamics, atmospheric conditions, and seawater salinity. The input
384 variables employed in our machine learning model encompass ERA5 data (sea level
385 pressure), NSIDC information (sea ice concentration and motion), as well as ETOPO1



386 dataset (bathymetric details). After filtering, we employ four machine learning
387 algorithms (Random Forest, K Nearest Neighbor, LightGBM, CatBoost) to train
388 salinity data obtained from EN4 and CTD. Utilizing multiple machine learning methods
389 can mitigate the impact of inherent flaws in a specific method on the results. During
390 data integration, varying weight combinations of variables greatly affect uncertainty;
391 therefore, we implement an uncertainty threshold to constrain appropriate weights.

392 We conducted an analysis to determine the significance of five input variables in
393 predicting salinity, which serves as a reliable indicator for identifying the key factors
394 influencing salinity changes. However, it is crucial to acknowledge that there might be
395 potential interactions among different variables. The importance of various factors
396 varies when predicting salinity in both EN4 and CTD datasets. Interestingly, both
397 datasets consistently highlight sea level pressure as the primary influential factor for
398 surface salinity prediction, while sea ice concentration emerges as the main determinant
399 when forecasting salinity at a depth of approximately 200m (corresponding to the
400 halocline base). The impact of sea ice movement on the surface is more significant than
401 that on the bottom of the halocline. The meridional ice speed is advantageous for
402 salinity prediction using CTD data, while the zonal flow speed is advantageous for
403 salinity prediction using EN4 data. However, the contribution of water depth factors
404 varies. CTD data indicates that water depth has a dominant influence on salinity
405 prediction in deep layers, whereas EN4 data shows the opposite trend. Salinity is closely
406 associated with freshwater distribution. The transport and accumulation of surface
407 freshwater are regulated by the sea level pressure field, and the melting of sea ice exerts
408 a greater impact on salinity compared to its movement affecting freshwater.



409



410 **Figure 11. Importance of different input variance.**

411 Accurate salinity product is crucial for understanding the dynamics of the Beaufort
412 Gyre and the redistribution of freshwater in the Beaufort Gyre in the western Arctic
413 Ocean. [Hall et al. \(2022\)](#) demonstrated that EN4 and ORAS5 salinity data can be
414 utilized for Arctic Ocean studies. However, when compared to EN4 and ORAS5,
415 salinity-derived freshwater content aligns more closely with BGEF estimates,
416 suggesting superior accuracy in FWC calculations. Furthermore, considering the
417 precision depth of halocline base, salinity products exhibit greater accuracy than EN4
418 and ORAS5. The findings from salinity product reveal a significant increase in
419 freshwater content throughout the upper 200m layer of the Beaufort Gyre during the
420 2000s; however, surface freshwater decreased while subsurface fresh water continued
421 to accumulate during the 2010s. It is likely that surface fresh water has been
422 redistributed towards Marklov Basin ([Bertosio et al., 2022](#)), potentially accumulating
423 in subsurface layers due to Ekman Pumping influences.

424 The salinity field of the Western Arctic Ocean is taken as an example to construct a
425 novel data mining method for polar sea areas, utilizing multiple machine learning
426 methods that integrate multiple data sources and incorporate physical processes. The
427 application potential of this method extends beyond the salinity field and includes other
428 related fields like hydrometeorology, sea ice thickness, polar biogeochemistry, among
429 others. It effectively utilizes multi-machine learning results for data evaluation and
430 integration.

431 **Author contributions.** LD provided scientific ideas, reviewed the paper and
432 contributed to the revising of figures and words of this paper; ST collected the datasets,
433 wrote the codes, analyzed the data, plotted the figures and wrote the paper. JL
434 contributed to the revising of figures and words of this paper

435 **Competing interests.** The contact author has declared that none of the authors has any
436 competing interests.

437 **Disclaimer.** Publisher's note: Copernicus Publications remains neutral with regard to
438 jurisdictional claims made in the text, published maps, institutional affiliations, or any
439 other geographical representation in this paper. While Copernicus Publications makes
440 every effort to include appropriate place names, the final responsibility lies with the
441 authors.

442 **Acknowledgements.** We collected a large number of public data, collected during



443 numerous expeditions and cruises over several decades. It is impossible to mention all
444 researchers who have contributed to these projects. We therefore thank all these people
445 for the gigantic effort, the huge amount of work and for making their data freely
446 available. This work was supported by the National Natural Science Foundation of
447 China (grant no. 42230405, 41976217, and 41576020) and the Global Change Research
448 Program of China (No. 2015CB953902).

449 6. Reference

- 450 Armitage, T. W., Bacon, S., Ridout, A. L., Petty, A. A., Wolbach, S. and Tsamados, M. (2017). Arctic
451 Ocean surface geostrophic circulation 2003–2014. *The Cryosphere*, 11(4), 1767–1780.
452 <https://doi.org/10.5194/tc-11-1767-2017>.
- 453 Behrendt, A., Sumata, H., Rabe, B., and Schauer, U.: UDASH – Unified Database for Arctic and
454 Subarctic Hydrography, *Earth Syst. Sci. Data*, 10, 1119–1138, <https://doi.org/10.5194/essd-10-1119-2018>, 2018.
- 456 Bertosio, C., Provost, C., Athanase, M., Sennéchaël, N., Garric, G., Lellouche, J. M., Bricaud, C.,
457 Kim, J.-H., Cho, K.-H. and Park, T. (2022). Changes in freshwater distribution and pathways in the
458 Arctic Ocean since 2007 in the Mercator Ocean global operational system. *Journal of Geophysical*
459 *Research: Oceans*, 127(6), e2021JC017701. <https://doi.org/10.1029/2021JC017701>
- 460 Carmack, E. C., Yamamoto-Kawai, M., Haine, T. W., Bacon, S., Bluhm, B. A., Lique, C., Melling,
461 H., Polyakov, I. V., Straneo, F., Timmermans, M. –L. and Williams, W. J. (2016). Freshwater and its
462 role in the Arctic Marine System: Sources, disposition, storage, export, and physical and
463 biogeochemical consequences in the Arctic and global oceans. *Journal of Geophysical Research:*
464 *Biogeosciences*, 121(3), 675–717. <https://doi.org/10.1002/2015JG003140>
- 465 Carmack, E., McLaughlin, F., Yamamoto-Kawai, M., Itoh, M., Shimada, K., Krishfield, R. and
466 Proshutinsky, A. (2008). Freshwater storage in the Northern Ocean and the special role of the
467 Beaufort Gyre. *Arctic–Subarctic ocean fluxes: defining the role of the northern seas in climate*, 145–
468 169. https://doi.org/10.1007/978-1-4020-6774-7_8.
- 469 Chen, L., Wang, Q., Zhu, G., Lin, X., Qiu, D., Jiao, Y., Lu, S., Li, R., Meng, G., and Wang, Y. (2024).
470 Dataset of stable isotopes of precipitation in the Eurasian continent. *Earth System Science Data*,
471 16(3), 1543–1557. <https://doi.org/10.5194/essd-16-1543-2024>
- 472 Cornish, S. B., Kostov, Y., Johnson, H. L. and Lique, C. (2020). Response of Arctic freshwater to
473 the Arctic oscillation in coupled climate models. *Journal of Climate*, 33(7), 2533–2555.
474 <https://doi.org/10.1175/JCLI-D-19-0685.1>
- 475 Giles, K. A., Laxon, S. W., Ridout, A. L., Wingham, D. J. and Bacon, S. (2012). Western Arctic
476 Ocean freshwater storage increased by wind-driven spin-up of the Beaufort Gyre. *Nature*
477 *Geoscience*, 5(3), 194–197. <https://doi.org/10.1038/ngeo1379>
- 478 Hall, S. B., Subrahmanyam, B., and Morison, J. H. (2021). Intercomparison of salinity products in
479 the Beaufort Gyre and Arctic Ocean. *Remote Sensing*, 14(1), 71.



-
- 480 <https://doi.org/10.3390/rs14010071>.
- 481 Holliday, N. P., Bersch, M., Berx, B., Chafik, L., Cunningham, S., Florindo-López, C., Hátún, H.,
482 Johns, W., Josey, S. A., Larsen, K. M. H., Mulet, S., Oltmanns, M., Reverdin, G., Rossby, T., Thierry,
483 V., Valdimarsson, H., and Yashayaev, I. (2020). Ocean circulation causes the largest freshening
484 event for 120 years in eastern subpolar North Atlantic. *Nature communications*, 11(1), 585.
485 <https://doi.org/10.1038/s41467-020-14474-y>.
- 486 Lin, P., Pickart, R. S., Heorton, H., Tsamados, M., Itoh, M. and Kikuchi, T. (2023). Recent state
487 transition of the Arctic Ocean's Beaufort Gyre. *Nature Geoscience*, 16(6), 485–491.
488 <https://doi.org/10.1038/s41561-023-01184-5>.
- 489 Manucharyan, G. E., Spall, M. A. and Thompson, A. F. (2016). A Theory of the Wind-Driven
490 Beaufort Gyre Variability. *Journal of Physical Oceanography*, 46(11), 3263–3278.
491 <https://doi.org/10.1175/jpo-d-16-0091.1>
- 492 Nguyen, A. T., Kwok, R. and Menemenlis, D. (2012). Source and Pathway of the Western Arctic
493 Upper Halocline in a Data-Constrained Coupled Ocean and Sea Ice Model. *Journal of Physical*
494 *Oceanography*, 42(5), 802–823. <https://doi.org/10.1175/jpo-d-11-040.1>
- 495 Proshutinsky, A., Krishfield, R., Timmermans, M. L., Toole, J., Carmack, E., McLaughlin, F.,
496 Williams, W. J., Zimmermann, S., Itoh, M. and Shimada, K. (2009). Beaufort Gyre freshwater
497 reservoir: State and variability from observations. *Journal of Geophysical Research: Oceans*,
498 114(C1). <https://doi.org/10.1029/2008JC005104>
- 499 Proshutinsky, A., Krishfield, R., Toole, J. M., Timmermans, M. L., Williams, W., Zimmermann, S.,
500 Yamamoto-Kawai, M., Armitage, T. W. K., Dukhovskoy, D., Golubeva, E., Manucharyan, G. E.,
501 Platov, G. , Watanabe, E. , Kikuchi, T. , Nishino, S., Itoh, M., Kang, S.-H., Cho, K.-H., Tateyama,
502 K. and Zhao, J. (2019). Analysis of the Beaufort Gyre freshwater content in 2003–2018. *Journal of*
503 *Geophysical Research: Oceans*, 124(12), 9658–9689. <https://doi.org/10.1029/2019JC015281>
- 504 Rabe, B., Karcher, M., Kauker, F., Schauer, U., Toole, J. M., Krishfield, R. A., Pisarev, S., Kikuchi,
505 T. and Su, J. (2014). Arctic Ocean basin liquid freshwater storage trend 1992–2012. *Geophysical*
506 *Research Letters*, 41(3), 961–968. Portico. <https://doi.org/10.1002/2013gl058121>.
- 507 Regan, H. C., Lique, C., and Armitage, T. W. (2019). The Beaufort Gyre extent, shape, and location
508 between 2003 and 2014 from satellite observations. *Journal of Geophysical Research: Oceans*,
509 124(2), 844–862. <https://doi.org/10.1029/2018jc014379>
- 510 Tao, S. and Du, L. (2024). Data mining-based machine learning methods for improving hydrological
511 data: a case study of salinity field in the Western Arctic Ocean [Data set]. Zenodo.
512 <https://doi.org/10.5281/zenodo.10990138>
- 513 Wang, Z., Wang, G., Guo, X., Bai, Y., Xu, Y. and Dai, M. (2022). Spatial reconstruction of long-
514 term (2003–2020) sea surface pCO₂ in the South China Sea using a machine learning based
515 regression method aided by empirical orthogonal function analysis. *Earth System Science Data*,
516 2023, 1–30. <https://doi.org/10.5194/essd-15-1711-2023>.
- 517 Zhang, J., Weijer, W., Steele, M., Cheng, W., Verma, T. and Veneziani, M. (2021). Labrador Sea
518 freshening linked to Beaufort Gyre freshwater release. *Nature communications*, 12(1), 1229.

<https://doi.org/10.5194/essd-2024-138>
Preprint. Discussion started: 3 May 2024
© Author(s) 2024. CC BY 4.0 License.



519 <https://doi.org/10.2172/1766967>