

1 **A historical nutrient dataset (1895–2024) for the North Pacific:**
2 **reconstructed from machine learning and hydrographic observations**

Chuanjun Du^{1*}, Naiwen Zheng¹, Shuh-Ji Kao¹, Minhan Dai², Zhimian Cao², Dalin
Shi², Qiancheng Li¹, Hao Wang¹, Xunlan Luo¹, and Xiaolin Li^{2*}

¹School of Marine Sciences, Hainan University, Haikou 570228, China

²State Key Laboratory of Marine Environmental Science, College of Ocean and Earth
Sciences, Xiamen University, Xiamen 361102, China

Manuscript resubmitted to *Earth System Science Data*

***Corresponding Authors:** Chuanjun Du, cjdu@hainanu.edu.cn; Xiaolin Li,
xlli@xmu.edu.cn

3

4 **Key points:**

- 5 ● Rigorous data quality control procedures were applied to clean nutrient and
6 hydrographic data collected from multiple sources in the North Pacific, following
7 state-of-the-art practices.
- 8 ● Three machine learning models demonstrated low errors across diverse validation
9 strategies.
- 10 ● We reconstructed a monumental database of ~473 million nutrient data points
11 across 1.92 million stations (1895–2024), expanding the number of nutrient data
12 points by a factor of 2,127–2,393 compared to original observations.

13

14

15 **Abstract**

16 Nutrients play a critical role in oceanic primary productivity and the biological pump.
17 However, compared to hydrographic parameters such as temperature and salinity,
18 nutrient observations are limited due to their labor-intensive and costly measurements.
19 Thus, nutrient observations are several orders of magnitude sparser than hydrographic
20 observations. In this study, we first established a rigorous data quality control procedure
21 to clean the hydrographic and nutrient (including NO_3^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$)
22 observations collected from World Ocean Database (WOD) and CLIVAR and Carbon
23 Hydrographic Data Office (CCHDO) in the North Pacific. Subsequently, the cleaned
24 and high-quality CCHDO dataset was used to train three machine learning models—
25 Random Forest, Light Gradient Boosting Machine (LightGBM), and Gaussian Process
26 Regression—to establish relationships between nutrient concentrations and key
27 variables, including space coordinates (longitude, latitude, and depth), time variables
28 (year and month), and water mass properties (indexed by potential temperature and
29 salinity). Validation shows that the reconstruction closely matches the observations,
30 with Root Mean Squared Errors (RMSEs) of <1.41 , <0.071 , <0.089 and <3.07
31 $\mu\text{mol kg}^{-1}$ for NO_3^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$, respectively. The validated models were
32 then applied to reconstruct nutrient concentrations from the hydrographic observations
33 in WOD, most of which lacked direct nutrient measurements. This resulted in ~ 473
34 million reconstructed nutrient data points across 1.92 million stations for each nutrient,
35 spanning from 1895 to 2024, representing a 2,127 to 2,393-fold increase compared to
36 the original nutrient observations in the North Pacific (197,539 to 222,234). This new
37 dataset will be valuable for studying nutrient transport and budgets, spinning up and
38 validating ocean biogeochemical models, assessing long-term nutrients and their
39 stoichiometric changes driven by anthropogenic forcing and climate change.~~variability~~
40 ~~under climate change and anthropogenic influences, and for providing transient~~
41 ~~boundary conditions in ocean biogeochemical models.~~ The dataset generated in this
42 study is openly available on Zenodo at <https://zenodo.org/records/17451417>.

43

44 **1 Introduction**

45 Bio-essential elements such as nitrogen, phosphorus, and silicon constitute the
46 fundamental material basis for marine ecosystems. Their concentrations govern
47 primary and new production (e.g., Browning et al., 2023; Lipschultz et al., 2002; Moore
48 et al., 2013) and subsequently regulate oceanic uptake of atmospheric CO₂ (Deutsch
49 and Weber, 2012; Sigman and Hain, 2012). However, traditional nutrient data collection
50 relies heavily on ship-based cruises and subsequent sample analysis, which are labor-
51 intensive, inefficient, and costly (Du et al., 2021). Consequently, compared to the
52 abundant hydrographic data collected from multiple platforms such as Conductivity-
53 Temperature-Depth (CTD) and the Array for Real-time Geostrophic Oceanography
54 (Argo) profilers, etc., nutrient observations are sparse in the ocean. These sparse
55 nutrient observations limit our understanding of both small-scale and long-term nutrient
56 variations and our comprehensive understanding of the mechanisms driving changes in
57 oceanic production and ecosystem dynamics (Bidigare et al., 2009; Yasunaka et al.,
58 2021; Karl et al., 2021).

59 To address this data sparsity, two main approaches have been commonly employed
60 to augment the spatiotemporal coverage of the observed nutrient data. The first is
61 objective analysis, which interpolates field measurements to generate broader spatial
62 coverage, as implemented in products such as the World Ocean Atlas (WOA) (e.g.,
63 Reagan et al., 2023; Lee et al., 2023). The second is data fusion, which establishes
64 statistical relationships between nutrients and environmental predictors such as
65 temperature (e.g., Kamykowski, 1987; Kamykowski et al., 2002; Kamykowski, 2008),
66 density (e.g., Dugdale et al., 1989; Switzer et al., 2003), oxygen, salinity, and
67 chlorophyll *a* (Goes et al., 1999; Palacios et al., 2013; Sarangi et al., 2011). Statistical
68 methods including cubic regression, multiple linear regression (Steinhoff et al., 2010;
69 Arteaga et al., 2015; Madani et al., 2024; Zhong et al., 2024), and generalized additive
70 models (Palacios et al., 2013) are frequently used in these efforts.

71 Recent studies have demonstrated the potential of machine learning for enhancing
72 the spatial and temporal coverage of nutrient data. For instance, Możejko and Gniot

73 (2008) used Artificial Neural Networks (ANNs) to model time series of total
74 phosphorous concentrations in the Odra River. Self-organizing maps (SOMs) were used
75 to estimate mixed layer nitrate and sea surface nutrients in the open ocean (Steinhoff et
76 al., 2010; Yasunaka et al., 2014). Liu et al. (2022) applied Support Vector Regression,
77 Random Forest Regression, and ANNs to reconstruct monthly surface nutrient
78 concentrations in the Yellow and Bohai Seas from 2003 to 2019. Their results revealed
79 pronounced seasonal and spatial variability in nutrient levels and underscored the
80 influence of environmental drivers such as sea surface temperature and salinity.
81 Similarly, Sundararaman and Shanmugam (2024) employed Gaussian Process
82 Regression (GPR) models to estimate global ocean surface macronutrient
83 concentrations using satellite-derived data, achieving high accuracy and demonstrating
84 their suitability for large-scale marine ecosystem monitoring. Yang et al. (2024)
85 employed a U-net and Earthformer to reconstruct the three-dimensional nitrate
86 distribution by integrating surface data including wind speed, sea surface temperature,
87 chlorophyll *a*, solar radiation, and precipitation in the Indian Ocean. These
88 advancements highlight the expanding role of machine learning in marine biochemical
89 data fusion and provide novel insights into nutrient dynamics and their ecological
90 impacts.

91 However, many existing approaches rely solely on mathematical extrapolation or
92 data fusion and often neglect the influence of physical seawater properties, such as
93 water mass characteristics. Using the relationship between nutrient concentration and
94 water masses (indexed by temperature and salinity), Du et al. (2021) successfully
95 predicted the nutrient concentrations in the South China Sea. However, the water
96 masses and their relationship with nutrients can also vary with space and time, which
97 should also be taken into consideration. In addition, most research has predominantly
98 focused on nutrient predictions at surface waters—driven by readily available remote-
99 sensing measurements of sea surface temperature and chlorophyll *a*—while subsurface
100 nutrient distributions remain poorly studied.

101 The North Pacific Ocean is one of the largest marine biomes in the global ocean (Karl
102 and Church, 2017), spanning a broader longitudinal range than the other oceans in the
103 world and a latitudinal range from tropical to subpolar regions. It includes a subtropical
104 gyre characterized by extremely low surface nutrient concentrations due to Ekman
105 convergence (e.g., Dave and Lozier, 2010; Browning et al., 2021; Dai et al., 2023), and
106 subpolar gyres in the north with elevated nutrient concentrations driven by Ekman
107 divergence. The atmospheric deposition (e.g., Martino et al., 2014; Qi et al., 2020), N₂-
108 fixation (e.g., Dai et al., 2023), and denitrification (Bonnet et al., 2017) are thought to
109 the main nutrient sources and sinks, which are decoupled in space and time in the North
110 pacific. It has been reported that the North Pacific Subtropical Gyre (NPSG) plays an
111 important role in fixed N inputs in summer, but also contributes disproportionately to
112 losses due to intense water-column denitrification in the eastern Pacific low-oxygen
113 zones (Eugster et al., 2012; Wang et al., 2019).

114 The North Pacific Ocean is influenced by multiple upwelling and current systems,
115 including the equatorial and California upwelling systems, North Equatorial Current,
116 Kuroshio Current, etc., which further change nutrient levels in these regions. In addition,
117 the North Pacific Ocean exhibits abundant mesoscale eddies (Chelton et al., 2007),
118 which play a critical role in redistributing nutrients and modulating biological activity
119 (e.g., Benitez-Nelson et al., 2007; Ascani et al., 2013; Barone et al., 2022). The
120 interaction of these multi-scale physical processes with biogeochemical processes
121 results in highly dynamic nutrient variability in the upper ocean. Therefore, high-
122 resolution and extensive nutrient datasets are essential to accurately resolve the nutrient
123 dynamics. Although the WOA (Reagan et al., 2023) serves as a primary nutrient
124 database and is widely used for boundary conditions in biogeochemical models, its
125 applicability is constrained by relatively coarse spatial resolution (currently 1°) and
126 climatological smoothing, which limit its ability to represent mesoscale and episodic
127 features or to capture long-term variations.

128 In the North Pacific, Yasunaka et al. (2014) used the SOMs technique to generate
129 monthly surface nutrient maps by integrating sea surface temperature, salinity,

130 chlorophyll *a*, and mixed layer depth. These maps revealed seasonal and interannual
131 variability in surface nutrient distributions in the northern North Pacific. To investigate
132 long-term changes, Yasunaka et al. (2016) applied Optimal Interpolation to analyze the
133 spatial and temporal evolution of surface nutrient concentrations. Lee et al. (2023)
134 provided spatiotemporally gridded nitrate and phosphate data in northwest Pacific from
135 1980 to 2019 using the spatiotemporal kriging technique. Wang et al. (2023) used the
136 deep neural network model to estimate nitrate concentrations in the upper northwestern
137 Pacific Ocean using temperature and salinity as the primary input parameters.

138 In this study, we first collected nutrient data from public databases and applied
139 rigorous quality control procedures. Using machine learning methods, we established
140 relationships between nutrient concentrations and water mass properties, spatial
141 coordinates, and temporal variables. We then evaluated the model performance through
142 a comprehensive error analysis. Finally, the validated models were applied to
143 reconstruct historical nutrient distributions across the North Pacific from 1895 to 2024.

144 **2 Data and Methods**

145 **2.1 Observation data**

146 Field observations were originally downloaded from the Climate and Ocean:
147 Variability, Predictability, and Change (CLIVAR) and Carbon Hydrographic Data
148 Office (CCHDO), which distributes vessel-based hydrographic data from programs
149 such as the World Ocean Circulation Experiment (WOCE), Joint Global Ocean Flux
150 Study (JGOFS), GO-SHIP, CLIVAR, and other repeat hydrography efforts
151 (<https://cchdo.ucsd.edu/>). In total, 631 cruises were collected in the North Pacific,
152 comprising 228,091, 197,617, 225,403, and 212,660 data points for $\text{NO}_3^- + \text{NO}_2^-$
153 (NO_x^-), NO_2^- , DIP, and $\text{Si}(\text{OH})_4$, respectively (Table 1). The dataset spans from 1973
154 to 2022 and was downloaded on October 1 2024; any updates made after this date were
155 not included in this study. The data cover a geographic range from 120.08°E to 95.17°W
156 and from 2.05°S to 60.25°N. The study domain was slightly extended into the South
157 Pacific to mitigate potential boundary effects during model development.

158

159
160
161
162
163

164 Table 1. Information on nutrients and their associated hydrographic data collected
165 from CLIVAR and Carbon Hydrographic Data Office (CCHDO) and the data
166 information after quality control (QC).

	Original data information		Data information after QC	
	Data	Stations	Data	Stations
Temperature	327792328502	1512745274	327688	15125
Salinity	328502341874	15274	328275	15269
NO _x ⁻	217725228094	9588	213962214943	90219120
NO ₂ ⁻	197617	8233	197539	8228
DIP	225403	9623	222234	9474
Si(OH) ₄	212660	8220	210447	8121

167 Hydrographic data for nutrient reconstruction were obtained from the World Ocean
168 Database (WOD; Mishonov et al., 2024), which compiles observations from various
169 platforms, including Autonomous Pinniped Bathythermograph (APB), Conductivity-
170 Temperature-Depth profiler (CTD), Drifting Buoy (DRB), Glider (GLD), Mechanical
171 Bathythermograph (MBT), Moored Buoy (MRB), Ocean Station Data (OSD), Profiling
172 Float (PFL), and Undulating Oceanographic Recorder (UOR). Since nutrient
173 reconstruction models rely on relationships with water masses, only samples containing
174 both temperature and salinity measurements were used; therefore, most APB
175 observations, which record only temperature, were excluded. Among these platforms,
176 CTD, OSD, and PFL provided the majority of usable data. Additionally, several
177 marginal seas—including the South China Sea, the Yellow Sea, the Sea of Japan, and
178 the Sea of Okhotsk—were excluded from this study because they are semi-enclosed
179 and strongly influenced by terrestrial inputs. The spatial domain was consistent with

180 that used for the CCHDO dataset, while the temporal coverage extended from 1875 to
 181 2024. In total, 577,215,683 data points from 2,284,448 stations across 40,113 original
 182 cruises were collected (Table 2)._

183 In addition, the OSD data before 1970 were extracted for nutrient validation in
 184 section 3.1. A total of 102,424, 125,142, 447,335, and 294,734 data points were
 185 collected for NO₃⁻, NO₂⁻, DIP, and Si(OH)₄, respectively.

186 Table 2. Information on hydrographic data collected from World Ocean Database, and
 187 the data information after quality control (QC). See main text for acronyms' full
 188 name.

Platform	Original data information			Data information after QC		
	Data	Stations	Cruises	Data	Stations	Cruises
APB	692302	46454	189	543714	37209	154
CTD	157914052	315177	8785	135584007	297036	8415
GLD	119302218	288840	384	69834989	285778	380
OSD	8885341	592225	21169	6942902	505780	17671
PFL	284781001	700798	9511	255423345	680531	9099
UOR	3373799	26699	7	3304158	25813	6
MRB	1459032	293734	65	1019565	88487	19
DRB	807938	20521	3	0	0	0
Total	577215683	2284448	40113	472652680	1920634	35744

189

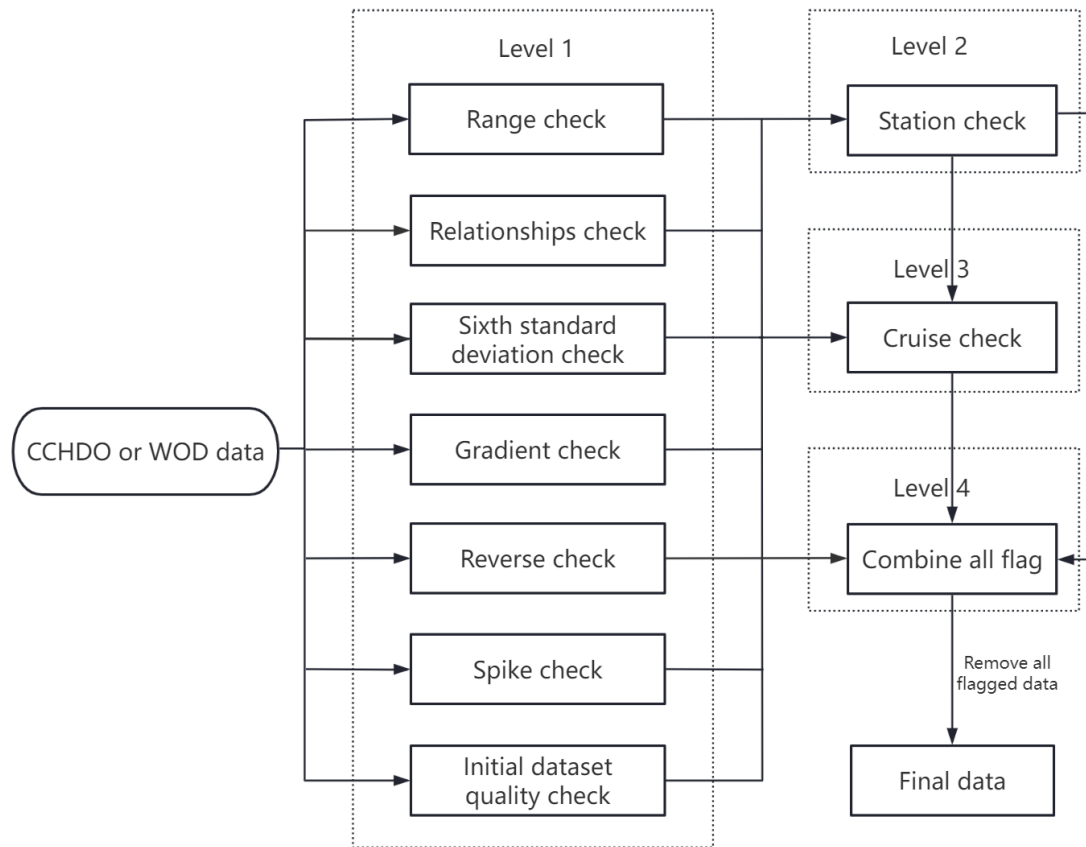
190

191 **2.2 Data quality control**

192 Given that the data were collected from multiple platforms using various methods
 193 over a long-time span and broad spatial range, quality control (QC) was essential (Du
 194 et al., 2021; Wang et al., 2025). Following the QC procedures developed by the World
 195 Ocean Database (WOD) (Garcia et al., 2024), we applied comprehensive QC protocols

196 (Fig. 21) to both CCHDO and WOD datasets, including hydrographic and nutrient
197 variables.

198 Four levels of QC were applied to identify and remove potentially erroneous or low-
199 quality records from the CCHDO and WOD datasets. The first level targeted individual
200 measurements, including several checks. (1) A range check was conducted by defining
201 depth-dependent acceptable value ranges for each parameter; data falling outside these
202 ranges were flagged as invalid. This check was applied to temperature, salinity, NO_x^- ,
203 NO_2^- , DIP, and $\text{Si}(\text{OH})_4$. Note that the NO_x^- denotes the sum concentration of NO_2^- and
204 NO_3^- . At stations lacking direct NO_x^- measurements, NO_x^- concentrations were derived
205 by summing discrete NO_2^- and NO_3^- observations. (2) An empirical relationship check
206 was performed to verify consistency among paired variables based on predefined
207 acceptable domains, including temperature–salinity, temperature– NO_x^- , temperature–
208 NO_2^- , temperature–DIP, temperature– $\text{Si}(\text{OH})_4$, salinity– NO_x^- , salinity– NO_2^- , salinity–
209 DIP, salinity– $\text{Si}(\text{OH})_4$, NO_x^- –DIP, and NO_x^- – $\text{Si}(\text{OH})_4$. (3) A six-standard-deviation
210 check was conducted by calculating the mean and standard deviation at each depth level;
211 values falling beyond six standard deviations were flagged as outliers. (4) A gradient
212 check assessed the vertical gradients of each parameter at each depth level across
213 stations; data showing abnormal gradients exceeding five standard deviations from the
214 mean were flagged as questionable. (5) A depth/potential density (σ_θ) inversion check
215 was applied to detect unrealistic reversals in parameters such as temperature and
216 nutrients, which typically exhibit monotonic relationships with depth or σ_θ in stratified
217 waters; measurements violating preset thresholds for depth–temperature, depth– NO_x^- ,
218 depth–DIP, depth– $\text{Si}(\text{OH})_4$, σ_θ –temperature, σ_θ – NO_x^- , σ_θ –DIP, and σ_θ – $\text{Si}(\text{OH})_4$ were
219 flagged. (6) A spike check was implemented to identify abrupt deviations (spikes)
220 between a measurement and its adjacent vertical neighbors; if the difference exceeded
221 a defined threshold, the data point was flagged as suspect. This check was applied to
222 temperature, NO_x^- , DIP, and $\text{Si}(\text{OH})_4$. (7) Only measurements with an original quality
223 flag of ‘good’ from CCHDO and WOD were retained, while those marked as
224 questionable or erroneous were flagged as outliers.



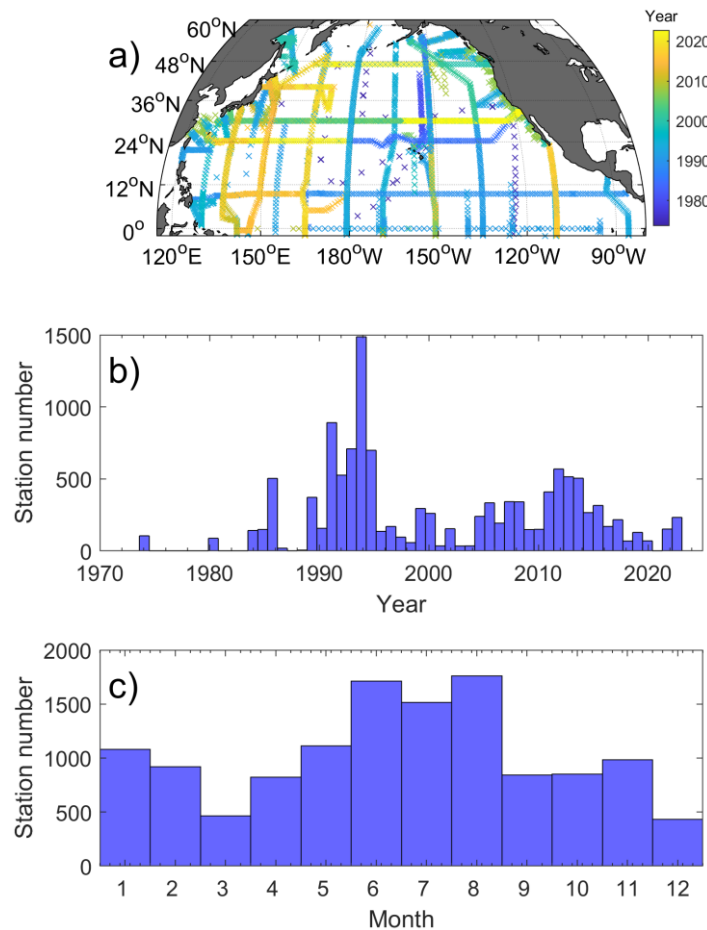
226

227 **Figure 1.** Data quality control procedures for temperature, salinity and nutrients
 228 collected from the CLIVAR and Carbon Hydrographic Data Office (CCHDO) and the
 229 World Ocean Database (WOD) datasets.

230

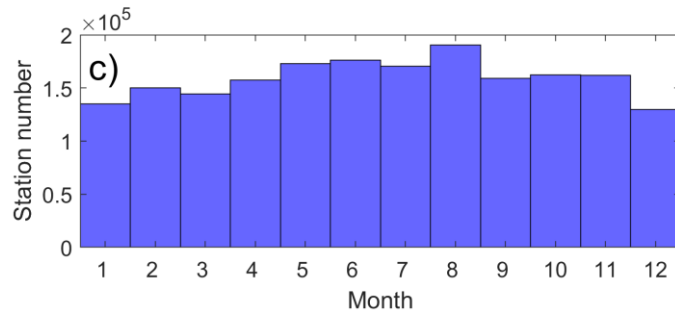
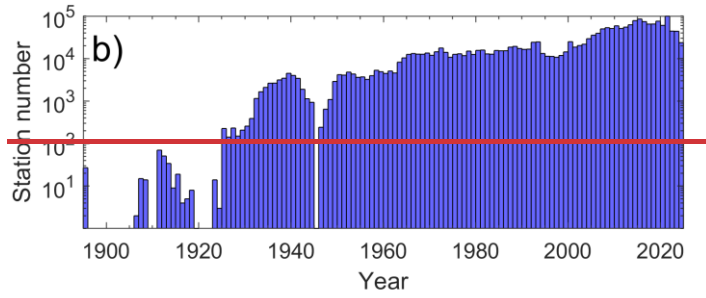
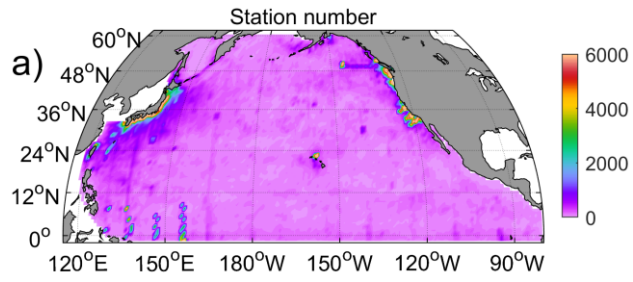
231 Building on the individual-level QC, we implemented additional QC at the station
 232 and cruise levels. At the station level, if a station profile contained more than 20%
 233 flagged data points, all data from that station were flagged as questionable. At the cruise
 234 level, if over 30% of a cruise's data were flagged, all data from that cruise were flagged.
 235 The final step integrated flags from all three levels (individual, station, and cruise), and
 236 any data flagged at any level were excluded. This hierarchical QC protocol effectively
 237 eliminates low-quality data. Although this approach may discard some high-quality
 238 measurements, the large volume of available data necessitates strict QC to ensure
 239 reliability.

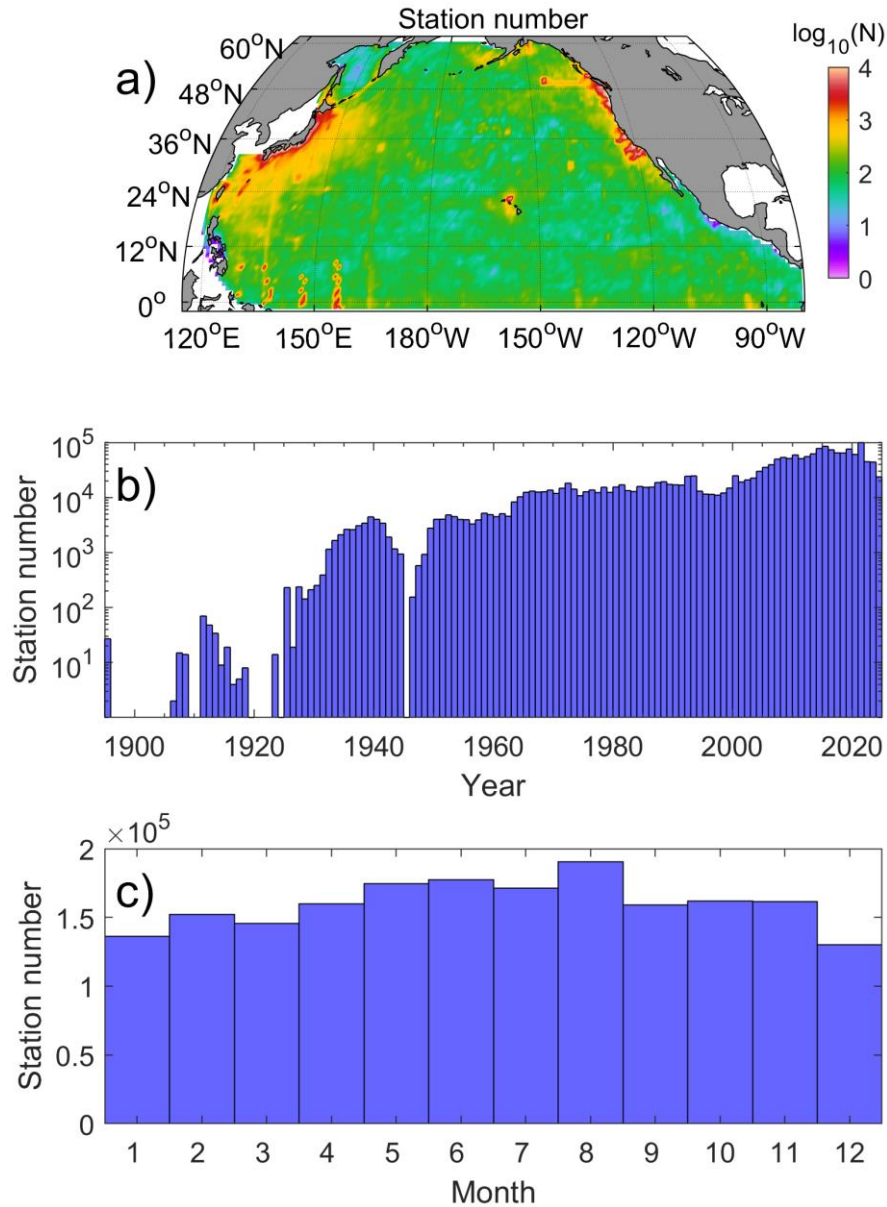
240 After quality control, the CCHDO dataset retained 214,943 (9,120), 197,539 (8,228),
 241 222,234 (9,457) and 210,447 (8,123) data points (stations), accounting for 94.2%
 242 (95.1%), 100.0% (99.9%), 98.6% (98.5%) and 99.0% (98.8%) of the original data
 243 points (stations) for NO_x^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$, respectively (Table 1). The retained
 244 stations cover nearly the entire North Pacific Ocean (Fig. 2a). ~~The retained data~~
 245 ~~spanninged~~ from 1972 to 2023. Most observations were collected after 1980, with a
 246 substantial increase after 1990 (Fig. 2b). Seasonally, the number of stations in June,
 247 July, and August was approximately three times greater than that in March and
 248 December (Fig. 2c).



249 **Figure 2.** Spatial and temporal distributions of NO_x^- (nitrate plus nitrite) after quality
 250 control in the North Pacific. a) Distribution of NO_x^- data locations, with points color-
 251 coded by year; b) station counts per year; c) station counts per month.
 252
 253

254 Following quality control, the final WOD dataset comprised 472,652,680
255 temperature and salinity data points from 1,920,634 stations across 35,744 cruises,
256 spanning 1895 to 2024. These represent 81.9% of the original observations, 84.1% of
257 the original stations, and 89.1% of the original cruises, respectively (Table 2). Spatially,
258 station counts per $1^\circ \times 1^\circ$ grid cell range from 1 to 31,851, with a mean of 249 stations
259 per cell (Fig. 3a). High sampling densities are found off eastern Japan and western
260 North America, resulting from high frequency observations from CTD and OSD
261 platforms, whereas elevated counts in the southwestern North Pacific primarily result
262 from MRB observations. Temporally, fewer than 300 stations per year were collected
263 before 1930. The annual number of stations exceeds 10,000 after 1964 and peaked at
264 approximately 100,000 in 2021 (Fig. 3b). Seasonally, station numbers are highest from
265 May to August (Fig. 3c). Overall, the collected WOD dataset provides 2127–2393 times
266 more observations and 202 times more station records than the CCHDO dataset.





268

269 **Figure 3.** Spatial and temporal distribution of the World Ocean Database (WOD) data
 270 after quality control. a) Station counts per $1^\circ \times 1^\circ$ grid cell; b) station counts per year;
 271 c) station counts per month.

272

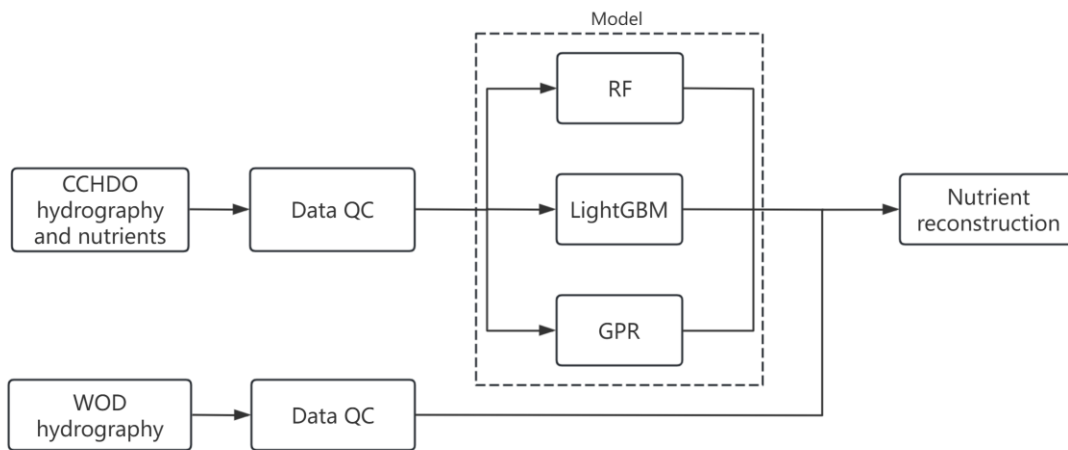
273 **2.3 Machine learning and nutrient reconstruction**

274 After rigorous data quality control, CCHDO data were used to train machine learning
 275 models. Three algorithms including Random Forest (RF), Light Gradient Boosting
 276 Machine (LightGBM), and Gaussian Process Regression (GPR) were applied to
 277 establish the relationship between environmental parameters and nutrient

278 concentrations. These methods are widely used in marine science (Hu et al., 2021;
279 Huang et al., 2022; Yu et al., 2022; Chen et al., 2023; Sundararaman and Shanmugam,
280 2024). The use of diverse models helps decrease algorithm selection bias. RF is an
281 ensemble technique based on bagging, which builds multiple independent decision
282 trees and aggregates their outputs by voting or averaging (Liaw and Wiener, 2002). Its
283 strengths include high predictive accuracy and reduced overfitting owing to the large
284 number of trees. RF has been applied to predict global primary production (Huang et
285 al., 2021), chlorophyll concentrations (Madani et al., 2024), nutrients (Chen et al., 2023;
286 Chen et al., 2024), dissolved iron (Huang et al., 2022), surface ocean $p\text{CO}_2$ (Chen et al.,
287 2019), and N_2 fixation rates (Yu et al., 2024).

288 LightGBM is an ensemble learning algorithm based on Gradient Boosting Decision
289 Trees (GBDT). Compared to standard GBDT, LightGBM employs a leaf-wise tree
290 growth strategy and a histogram-based binning technique to improve predictive
291 accuracy and computational efficiency (Ke et al., 2017). It has been successfully
292 applied to predict water levels (Gan et al., 2021), salinity (Dong et al., 2022; Wang et
293 al., 2022), and chlorophyll a concentration (Su et al., 2021). GPR is a non-parametric
294 Bayesian approach that infers relationships by defining a prior distribution over
295 functions via kernel-based covariance matrices, rather than estimating fixed
296 coefficients. This flexibility allows GPR to capture complex, nonlinear input–output
297 relationships and to quantify prediction uncertainty. GPR has been used in
298 oceanography to estimate global dissolved oxygen and nutrient concentrations
299 (Sundararaman and Shanmugam, 2024).

300

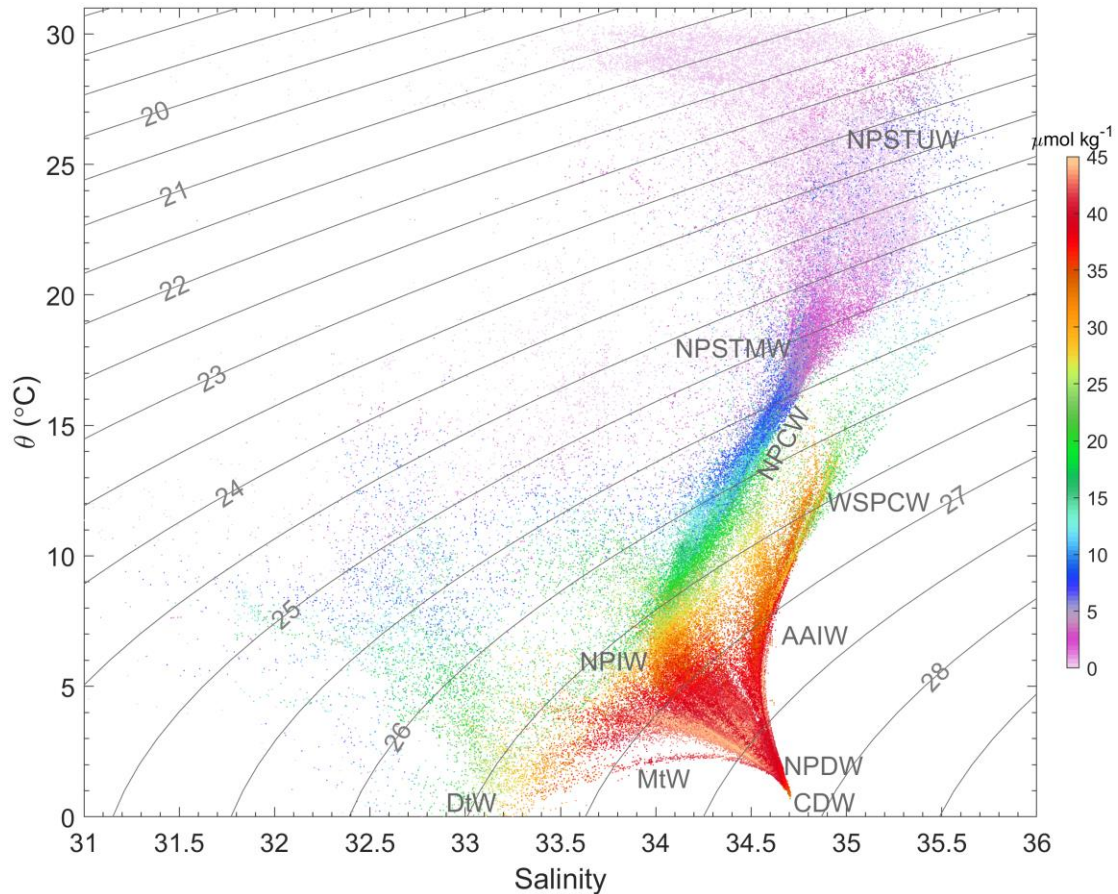


301

302 **Figure 4.** Flowchart of the machine learning framework and its application to WOD
303 hydrographic data for nutrient reconstruction.

304

305 In this study, we used spatial coordinates (longitude, latitude, depth), temporal
306 variables (month and year), and water mass properties (represented by potential
307 temperature and salinity) as environmental predictors of nutrient concentrations. The
308 time predictors used month and year with decimals to capture seasonal, interannual,
309 and long-term variability. The North Pacific contains distinct water masses, including
310 North Pacific Subtropical Water, North Pacific Intermediate Water, Antarctic
311 Intermediate Water, Western South Pacific Central Water, North Pacific Deep Water,
312 and Pacific Deep Water, as well as Circumpolar Deep Water (e.g., Talley et al., 2011;
313 Fuhr et al., 2021). These water masses mix to form different water types associated with
314 distinct nutrient concentrations (Fig. 5). Water types have been found to be an important
315 parameter to reconstruct nutrient concentrations in the South China Sea (Du et al., 2021).
316 Thus, potential temperature and salinity serve as proxies for water mass identification.



317

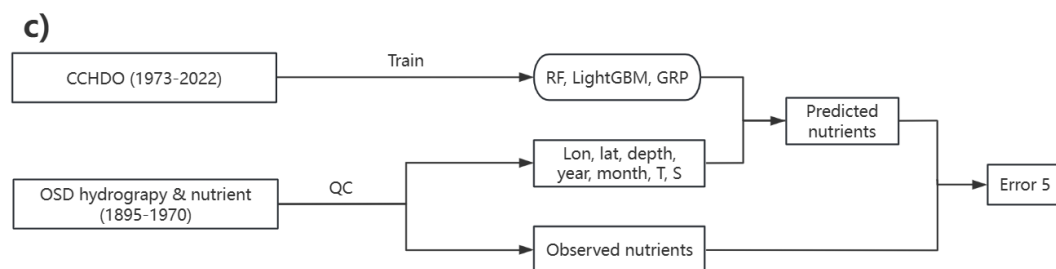
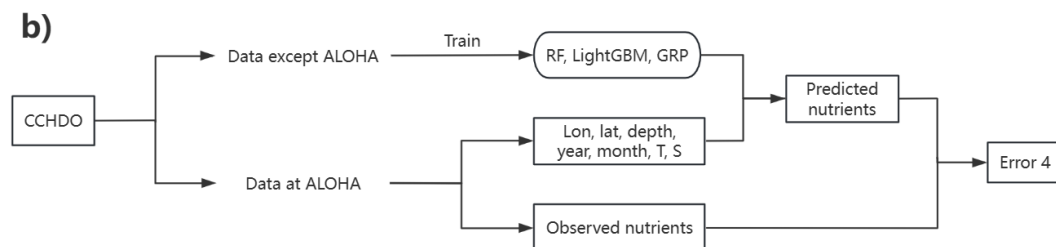
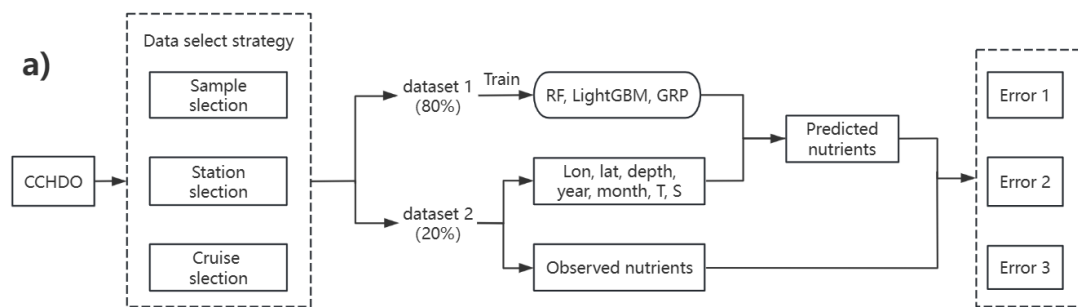
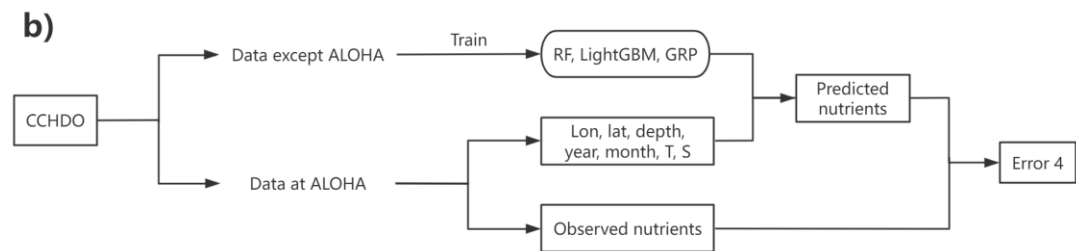
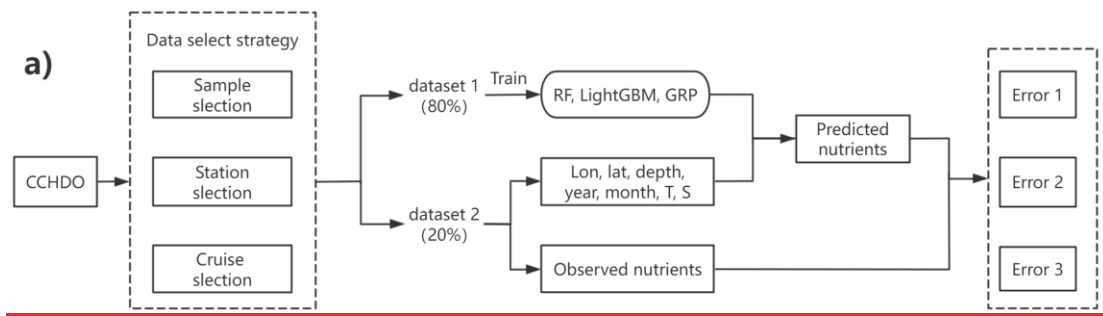
318 **Figure 5.** The water masses (indicated by salinity and potential temperature (θ)) and
 319 NO_x^- ($\text{NO}_3^- + \text{NO}_2^-$; color shading) relationships in the North Pacific. The temperature
 320 and salinity data were collected from the CCHDO dataset. The gray contour lines and
 321 number denote the potential density anomaly. The typical water masses are shown as
 322 follows: North Pacific Central Water (NPCW), North Pacific Subtropical Underwater
 323 (NPSTUW), North Pacific Subtropical Mode Water (NPSTMW), North Pacific
 324 Intermediate Water (NPIW), Dichothermal Water (DtW), Mesothermal Water (MtW),
 325 Antarctic Intermediate Water (AAIW), Western South Pacific Central Water (WSPCW),
 326 Pacific Deep Water (PDW), and Circumpolar Deep Water (CDW). The water masses
 327 and their acronyms are follow the classifications in Talley et al. (2011) and Fuhr et al.
 328 (2021).

329

330 **3 Results**

331 **3.1 Error estimation**

332 Leave-one-out cross-validation was primarily used to quantify model reconstruction
333 errors. The CCHDO dataset was divided into training and testing subsets for model
334 development and performance evaluation, respectively. To assess how data partitioning
335 affects error metrics, we implemented four validation methods based on different data-
336 selection strategies (Fig. 6a). The first three methods involved partitioning the CCHDO
337 dataset into training (80%) and testing (20%) subsets. These methods employed three
338 data selection strategies: (1) sample-random, by withholding 20% of individual samples;
339 (2) station-random, by withholding 20% of stations; and (3) cruise-random, by
340 withholding 20% of cruises. Predictions for the held-out subsets, generated using their
341 respective spatial, temporal, and water mass property data, were compared against the
342 actual withheld nutrient measurements to calculate error metrics. These partitioning
343 strategies were designed to evaluate potential errors under the sparse and non-uniform
344 spatiotemporal distribution of observations: Error 1 represented an optimistic estimate
345 (validation data are likely colocated with training data in space and time), Error 3
346 represented a conservative, generalized scenario (validation data are independent of
347 training data), Error 2 provided an intermediate estimate (validation data may share
348 spatial/temporal context with training data within the same cruise). The choice of error
349 metric (Error 1, 2, or 3) should be guided by the degree of extrapolation in the intended
350 application relative to the training data's spatiotemporal distribution.



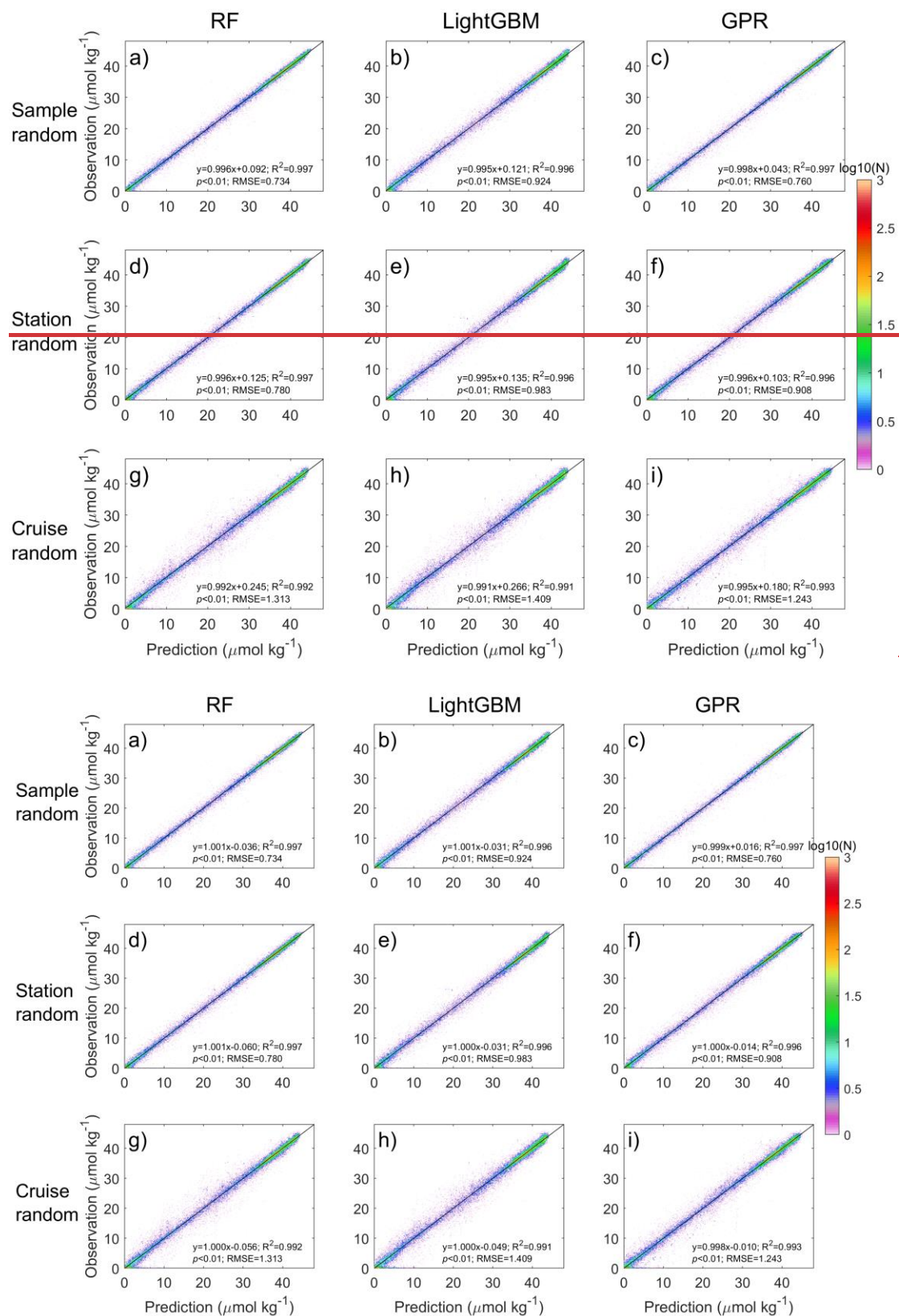
351

352

353 **Figure 6.** Schematic of the error estimation procedure. a) Error estimation based on
 354 three types of data selection strategy; b) assessing temporal error evolution by
 355 excluding the data at Station ALOHA; -c) examining the models' reconstruction error

356 using the hydrographic and nutrient data before 1970. The T and S denote the potential
357 temperature and salinity, respectively.

358
359 The validation results for reconstructed NO_x^- versus observations under the first three
360 data-selection strategies are shown in Fig. 7. RF and GPR exhibited nearly identical
361 performance, with regression slopes of 0.992–0.998, $R^2 > 0.992$, and Root Mean
362 Squared Errors (RMSE_s) between 0.734 and 1.313 $\mu\text{mol kg}^{-1}$ (Fig. 7a, c, d, f, g, i).
363 LightGBM showed slightly lower accuracy (slope: 0.991–0.995; R^2 : 0.991–0.996;
364 RMSE_s: 0.780–1.419 $\mu\text{mol kg}^{-1}$) (Fig. 7b, e, h). Across different data-selection
365 strategies, sample-random (Error 1) yielded the lowest errors (RMSE_s: 0.734–0.983
366 $\mu\text{mol kg}^{-1}$) (Fig. 7a–c), station-random (Error 2) was intermediate (RMSE_s: 0.908–
367 1.313 $\mu\text{mol kg}^{-1}$) (Fig. 7d–f), and cruise-random (Error 3) produced the highest errors
368 (RMSE_s: 1.243–1.424 $\mu\text{mol kg}^{-1}$) (Fig. 7; Table 3). This gradient in error estimates
369 underscores the necessity of employing different data-selection strategies for a
370 comprehensive error assessment. The high slopes and R^2 values (>0.99) achieved across
371 all algorithms and data-selection strategies confirmed the robustness of the nutrient
372 reconstructions.



373

374

375 **Figure 7.** Validating the reconstructed NO_x^- concentrations using leave-one-out cross-
 376 validation with different data selection strategies and machine learning methods. Plots
 377 shown in row 1 correspond to the sample random strategy (a-c), row 2 correspond to

378 the station random strategy (d-e), and row 3 correspond to the cruise random
 379 strategy (g-i). Plots shown in column 1 correspond to the Random Forest (RF; a, d, and
 380 g), column 2 correspond to the LightGBM (b, e, and h), and column 3 correspond to
 381 the Gaussian Process Regression (GPR; c, f, and i). The black lines and text show the
 382 fitted linear regressions, regression equations, coefficient of determination (R^2), p
 383 values, and Root Mean Squared Errors (RMSEs). The color represents the data density
 384 (N, number of observations). Note that the logarithmic scale of N is applied.

385

386 Reconstruction errors for NO_2^- , DIP, and $\text{Si}(\text{OH})_4$ are summarized in Figs. S1–S3
 387 and Table 3. Across methods, the RMSEs-values were below $0.079 \mu\text{mol kg}^{-1}$ for NO_2^- ,
 388 $0.089 \mu\text{mol kg}^{-1}$ for DIP, and $3.07 \mu\text{mol kg}^{-1}$ for $\text{Si}(\text{OH})_4$. DIP and $\text{Si}(\text{OH})_4$ exhibited
 389 similar error trends: RMSEs increased from sample-random to station-random to
 390 cruise-random selection. In contrast, NO_2^- reconstruction exhibited lower accuracy than
 391 NO_x^- , DIP, and $\text{Si}(\text{OH})_4$, with regression slopes of 0.48–0.68 and R^2 values of 0.32–
 392 0.72. RF and LightGBM outperform GPR for NO_2^- . The poorer NO_2^- performance
 393 likely reflects its generally low concentrations (mostly $<0.5 \mu\text{mol kg}^{-1}$) and high
 394 biological variability. Thus, we highlight NO_2^- as a high-uncertainty reconstruction.

395 Table 3 The Root Mean Squared Errors of nutrient reconstruction from different error
 396 evaluation strategies (unit: $\mu\text{mol kg}^{-1}$).

Data selection strategy	NO_x^-			NO_2^-			DIP			$\text{Si}(\text{OH})_4$		
	RF	Light GBM	GPR	RF	Light GBM	GPR	RF	Light GBM	GPR	RF	Light GBM	GPR
Sample random	0.724	0.924	0.760	0.049	0.054	0.079	0.056	0.070	0.055	1.90	2.30	1.53
Station random	0.780	0.983	0.908	0.065	0.068	0.072	0.058	0.071	0.065	2.07	2.45	2.20
Cruise random	1.313	1.409	1.243	0.054	0.057	0.071	0.080	0.089	0.084	2.79	3.07	2.94
ALOHA validation	0.701	0.842	0.674	—	—	—	0.066	0.079	0.064	2.13	2.48	2.32

397

398 Understanding the spatiotemporal structure of reconstruction errors is also important
 399 for assessing the models' reconstruction applicability. As shown in Figs. S4-S7, the
 400 reconstruction errors of NO_3^- , DIP, and $\text{Si}(\text{OH})_4$ are generally small in the surface layer,
 401 increase with depth to maxima at the nutricline, and then decrease to low values in deep

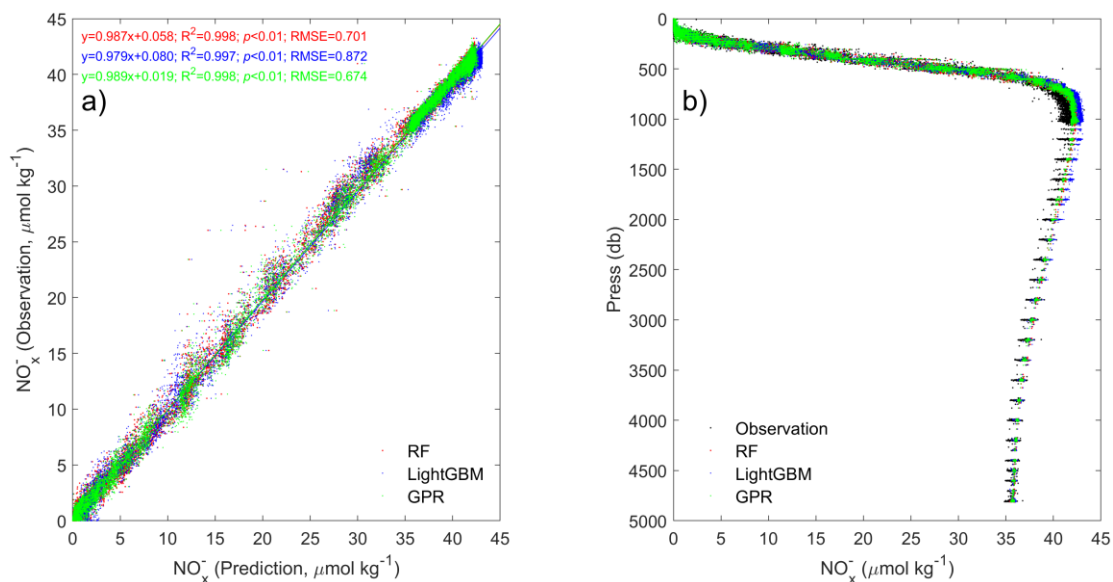
402 layers. However, the random errors associated with individual cruise observations for
403 Si(OH)₄ display no evident vertical pattern. Horizontally, we paid particular attention
404 to surface waters due to their greatest concentration gradients. The horizontal
405 distribution shows that the errors are small in the western NPSG (a nutrient-depleted
406 region) but are large in the subarctic gyre and close to the equatorial regions (nutrient-
407 replete regions; Figs. S8-S11). Here, we particularly examined the nutrient
408 reconstruction errors in the oligotrophic NPSG. The oligotrophic regimes are defined
409 as regions where NO₃⁻, NO₂⁻, DIP, and Si(OH)₄ concentrations are <0.2, <0.2, <0.2,
410 and <5.0 μmol kg⁻¹, respectively. As shown in Table 4, the reconstruction errors in these
411 regimes are <0.574, <0.056, <0.084, and <1.88 μmol kg⁻¹ for NO₃⁻, NO₂⁻, DIP, and
412 Si(OH)₄, respectively, which are evidently lower than the overall RMSEs for the entire
413 North Pacific (Table 3). Among these models, the RF generally performs the best
414 compared to the others. This confirms that absolute errors decrease in oligotrophic
415 regimes. Since the number of summer observations is up to three times greater than that
416 in winter and spring, we further examined the seasonal variation of errors. Overall, no
417 evident seasonal variations are displayed. Only in the case of random cruise selection
418 was the NO₃⁻ error shown to be greater in spring (March to May) than in other seasons
419 (Fig. S12). For other cases and nutrients, seasonal variation in error was not evident.
420 On a decadal timescale, the reconstruction errors display a slight decreasing trend,
421 particularly for DIP, from 1973 to 2020 (Fig. S13), implying that the errors might be
422 smaller in recent decades than in previous ones.

423 Table 4 The Root Mean Squared Errors of nutrient reconstruction from different error
424 evaluation strategies in surface oligotrophic regimes (unit: μmol kg⁻¹).

<u>Data</u> <u>selection</u> <u>strategy</u>	<u>NO_x⁻</u>			<u>NO₂⁻</u>			<u>DIP</u>			<u>Si(OH)₄</u>		
	<u>RF</u>	<u>Light</u> <u>GBM</u>	<u>GPR</u>	<u>RF</u>	<u>Light</u> <u>GBM</u>	<u>GPR</u>	<u>RF</u>	<u>Light</u> <u>GBM</u>	<u>GPR</u>	<u>RF</u>	<u>Light</u> <u>GBM</u>	<u>GP</u> <u>R</u>
<u>Sample</u> <u>random</u>	<u>0.290</u>	<u>0.567</u>	<u>0.444</u>	<u>0.018</u>	<u>0.035</u>	<u>0.048</u>	<u>0.028</u>	<u>0.042</u>	<u>0.039</u>	<u>1.19</u>	<u>0.90</u>	<u>1.30</u>
<u>Station</u> <u>random</u>	<u>0.303</u>	<u>0.457</u>	<u>0.474</u>	<u>0.030</u>	<u>0.030</u>	<u>0.043</u>	<u>0.036</u>	<u>0.045</u>	<u>0.043</u>	<u>1.24</u>	<u>1.51</u>	<u>1.51</u>
<u>Cruise</u> <u>random</u>	<u>0.378</u>	<u>0.457</u>	<u>0.574</u>	<u>0.030</u>	<u>0.029</u>	<u>0.056</u>	<u>0.075</u>	<u>0.077</u>	<u>0.084</u>	<u>1.85</u>	<u>1.88</u>	<u>1.75</u>

425
426

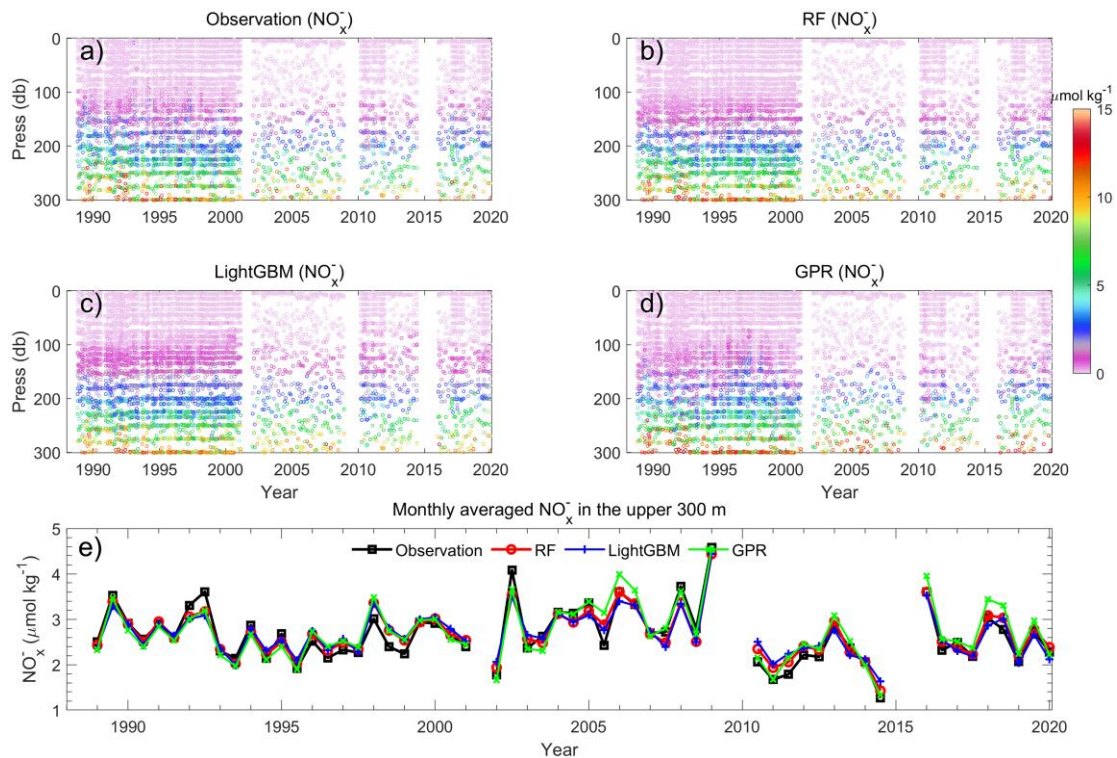
427 A fourth validation step assessed the model's temporal performance at Station
 428 ALOHA (Error 4; Fig. 6b). To test this, we withheld all observations from ALOHA
 429 (which, since 1988, represent 8.52%, 8.45%, and 8.11% of the total $\text{Si}(\text{OH})_4$, NO_x^- , and
 430 DIP records, respectively) from model training. We then reconstructed nutrient
 431 concentrations using space, time, and water-type predictors at Station ALOHA. NO_2^-
 432 was excluded due to insufficient observations. For NO_x^- , the regression slopes between
 433 reconstruction and observations were 0.99, 0.98, and 0.99, with RMSEs of 0.701, 0.842,
 434 and 0.674 $\mu\text{mol kg}^{-1}$ for RF, LightGBM, and GPR, respectively; R^2 values exceeded
 435 0.997 for all models (Fig. 8a). RF and GPR slightly outperformed LightGBM. All
 436 models accurately reproduced the NO_x^- profiles (Fig. 8b). The reconstruction errors for
 437 DIP were 0.066, 0.079, and 0.064 $\mu\text{mol kg}^{-1}$ for RF, LightGBM, and GPR, respectively.
 438 The corresponding errors for $\text{Si}(\text{OH})_4$ were 2.13, 2.48, and 2.32 $\mu\text{mol kg}^{-1}$ (Table 3,
 439 Figs. S14–S6S15).



440
 441 **Figure 8.** Validating the reconstructed nutrient concentrations at Station ALOHA. a)
 442 Reconstructed $\text{NO}_3^- + \text{NO}_2^-$ (NO_x^-) vs. observations: Random Forest (RF; red dots),
 443 LightGBM (blue dots), and Gaussian Process Regression (GPR; green dots). b) Profiles
 444 of observed (black dots) and reconstructed NO_x^- from RF (red dots), LightGBM (blue
 445 dots), and GPR (green dots).

446
 447 Since the variations of nutrients primarily occur in the upper water column, we
 448 focused on the nutrient reconstruction in the upper 300 m at Station ALOHA. Overall,
 449 the models reproduced the profiles of NO_x^- from 1988 to 2021 well (Fig. 9a-d). The

450 reconstruction errors were low at the surface and increased with depth, with most of the
 451 values $< 3.0 \mu\text{mol kg}^{-1}$ (Fig. S16a-d). To evaluate models' ability to reconstruct nutrient
 452 variations in time, the nutrient concentrations were averaged monthly over the upper
 453 300 m. As compared to observations, RF, LightGBM, and GPR all well reconstructed
 454 the interannual variations of NO_x^- , with most of the absolute errors $< 0.5 \mu\text{mol kg}^{-1}$
 455 (Figs. 9e and S16e) at Station ALOHA. Similarly, the validation of DIP and $\text{Si}(\text{OH})_4$
 456 are shown in Figs. S17-S20 at Station ALOHA (Figs. 9e, S6, and S7).



457
 458 **Figure 9.** Temporal variations of NO_x^- concentrations in the upper 300 m at Station
 459 ALOHA from 1988 to 2021 for observed (a) and reconstructed NO_x^- by Random Forest
 460 (RF; b), LightGBM (c), and Gaussian Process Regression (GPR; d). (e) Time series of
 461 monthly averaged NO_x^- concentrations in the upper 300 m from observations, and
 462 reconstructions by RF, LightGBM, and GPR.

463
 464 A fifth validation step evaluates the models' reconstruction for the period before 1970
 465 (Error 5; Fig. 6c). This is necessary because the training data (CCHDO) spans 1973–
 466 2022, while the reconstructions are extrapolated back to 1895. We argue that this
 467 extrapolation should be reasonable because the variations of temperature-salinity-
 468 nutrient relationships in the ocean's interior might be small over the past century,

469 providing a basis for temporal extrapolation. First, the residence time of nitrogen in
470 deep and intermediate waters can be up to 2000 years in the North Pacific. Consequently,
471 the imprint of centennial-scale change on nutrient inventories is attenuated. Second, the
472 long-term variations of nutrient concentrations are not evident within our core training
473 period (1973–2022; Figs. 9e and 17). Finally, the mean nutrient profiles derived from
474 the 1920-1970 and 1973-2022 periods are not evidently different in the central North
475 Pacific (Fig. S21). Therefore, while the North Pacific may experience long-term
476 variability, it might be masked by the reconstruction error, and the use of hydrographic
477 properties as predictors for nutrients is justified for historical reconstructions.

478 However, when assessing the reconstruction errors before 1970, we first consider
479 data quality issues. Prior to the standardization of modern oceanographic methods,
480 nutrient measurements—particularly from earlier decades—were subject to greater
481 analytical errors, inconsistent sampling protocols, and varied determination techniques.
482 The data quality concern is evident in the sporadic and sometimes physically
483 implausible deep nutrient profiles found in WOD for that era (Fig. S22). This is also
484 the primary reason that nutrient data pre-1973 collected from sources like the OSD from
485 WOD were not incorporate into model training. To evaluate data quality in earlier
486 decades, we selected five specific years with more abundant observations: 1929, 1947,
487 1953, 1958, and 1966 (Fig. S23). After applying the same quality-control criteria
488 outlined in Section 3.1, we used the historical hydrographic data (temperature and
489 salinity) from those years to predict nutrient concentrations. A total of 52,277, 119,137,
490 284,472, and 193,339 data points were collected for NO_3^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$,
491 respectively, after QC. The comparison between these predictions and the quality-
492 controlled observations yields the prediction errors for the pre-1970 period (Fig. 6c).
493 The RMSEs from different models suggested values <5.7 , <0.40 , and $<22.9 \mu\text{mol kg}^{-1}$
494 for NO_3^- , DIP, and $\text{Si}(\text{OH})_4$, respectively (Figs. S24–S26), which are much larger than
495 the corresponding errors for the period after 1970. We recommend that these values be
496 considered a conservative estimate of the upper error bound, as they incorporate both
497 nutrient observations and prediction errors. In addition, the hydrographic data are also

498 less reliable in the earlier period. Thus, we acknowledge that reconstruction errors are
499 likely higher for the pre-1973 period, and the error estimated here should be considered
500 as a "best estimate" with quantified uncertainties, and encourage users to consider these
501 error bounds when applying the dataset to early twentieth-century conditions.

503 **3.2 Reconstructed ~~_~~ nutrients and their distributions**

504 The final reconstructed nutrient dataset aligns with the spatiotemporal coverage of
505 the quality-controlled WOD hydrographic dataset, comprising 472,652,680 data points
506 for each nutrient (NO_x^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$) from 1,920,634 stations across 35,744
507 cruises, spanning from 1895 to 2024 (Table 2). Most data points are located above 2,000
508 m, with fewer observations at greater depths due to **observational**
509 **hydrographic** platform limitations. ~~Since the distribution patterns of NO_x^- , DIP, and~~
510 ~~$\text{Si}(\text{OH})_4$ are consistent across the different methods (Figs. 10–13, S8–S16), we focus~~
511 ~~on the reconstructed NO_x^- from RF model in this section unless stated otherwise.~~

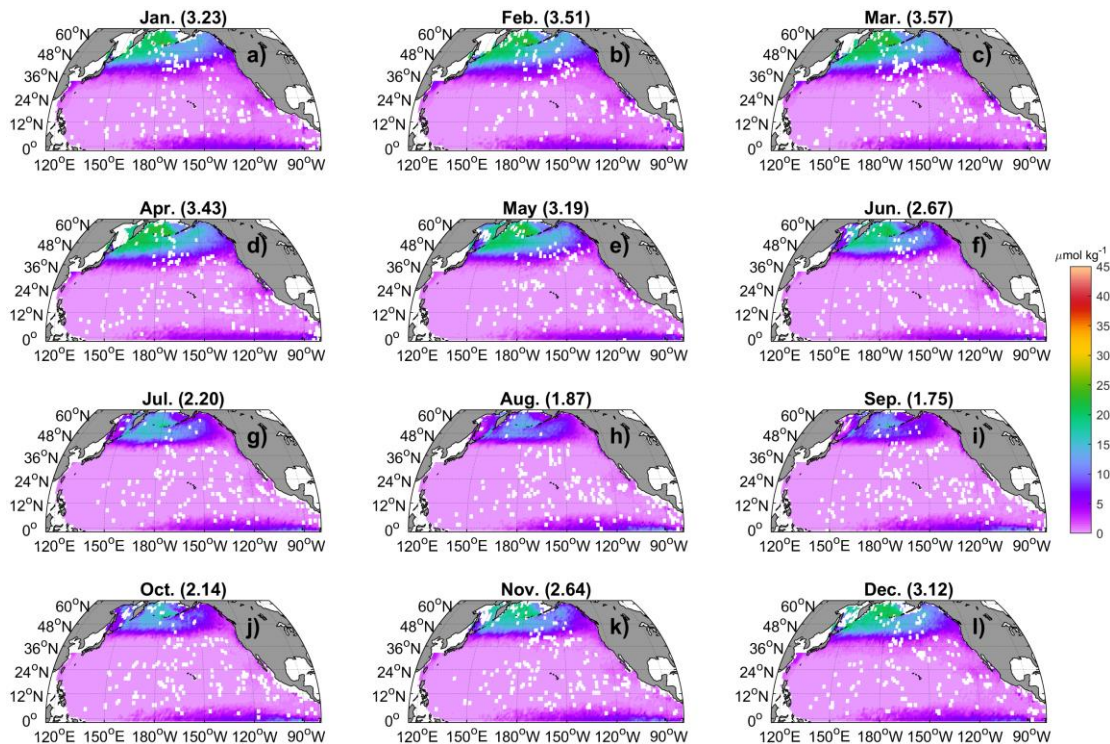
512 It is important to clarify the nature of the reconstructed dataset, which is
513 fundamentally different from gridded products. This product provides nutrient
514 concentrations linked to each hydrographic observations: nutrient values are
515 reconstructed precisely at the locations, depths, and times of original hydrographic
516 observations (sourced from WOD) where direct nutrient measurements might be
517 unavailable or of poor quality. This approach yields a point-wise dataset that aligns with
518 the original hydrographic observations, rather than a spatially or temporally
519 interpolated field—an important distinction for users interpreting and applying the data.

520 **3.3 Climatology of nutrient distributions**

521 ~~It should be noted that our climatology is derived from the mean of existing data,~~
522 ~~which heavily relies on the spatiotemporal distribution of those data and may not~~
523 ~~represent the true climatological mean. Since the distribution patterns of NO_x^- , DIP,~~
524 ~~and $\text{Si}(\text{OH})_4$ are consistent across the different methods (Figs. 10–13, S8–S16), we~~
525 ~~focus on the reconstructed NO_x^- from RF model in this section unless stated otherwise.~~

526 To evaluate the reliability of our product, we binned and averaged the predicted
527 nutrients within 1°×1° grid cells for each month to produce a monthly climatology. This
528 climatology represents a mean field that depends heavily on the spatiotemporal
529 distribution of the underlying data and may be influenced by uneven data sampling.
530 This reconstructed climatology was compared with the World Ocean Atlas 2023
531 (WOA23), which is derived from quality-controlled and objectively analyzed
532 observational data. Since the large-scale patterns of NO₃⁻, DIP, and Si(OH)₄ are similar
533 among different models (Figs. 10–13, S27–S36), we focus on NO₃⁻ reconstructed by
534 the RF model in this section unless stated otherwise.

535 Figs. 10–13 present the monthly climatology of NO_x⁻ at 5 m, 100 m, 500 m, and
536 1,000 m in the North Pacific. At 5 m, the reconstructed NO_x⁻ accurately captures the
537 established spatial patterns, with elevated concentrations in the subpolar gyre, Bering
538 Sea, and equatorial regions, and depleted concentrations in the ~~North-Pacific~~
539 ~~Subtropical Gyre (NPSG (Fig. 10))~~. Seasonally, the basin-averaged surface NO_x⁻
540 concentrations display the highest value of 3.50 μmol kg⁻¹ in March, in contrast to the
541 lowest value of 1.82 μmol kg⁻¹ in September. These results agree with Yasunaka et al.
542 (2014, 2021), who, using extensive surface nutrient observations (up to 14,000 for
543 nitrate) in the North Pacific, reported similar spatial and seasonal patterns.



544

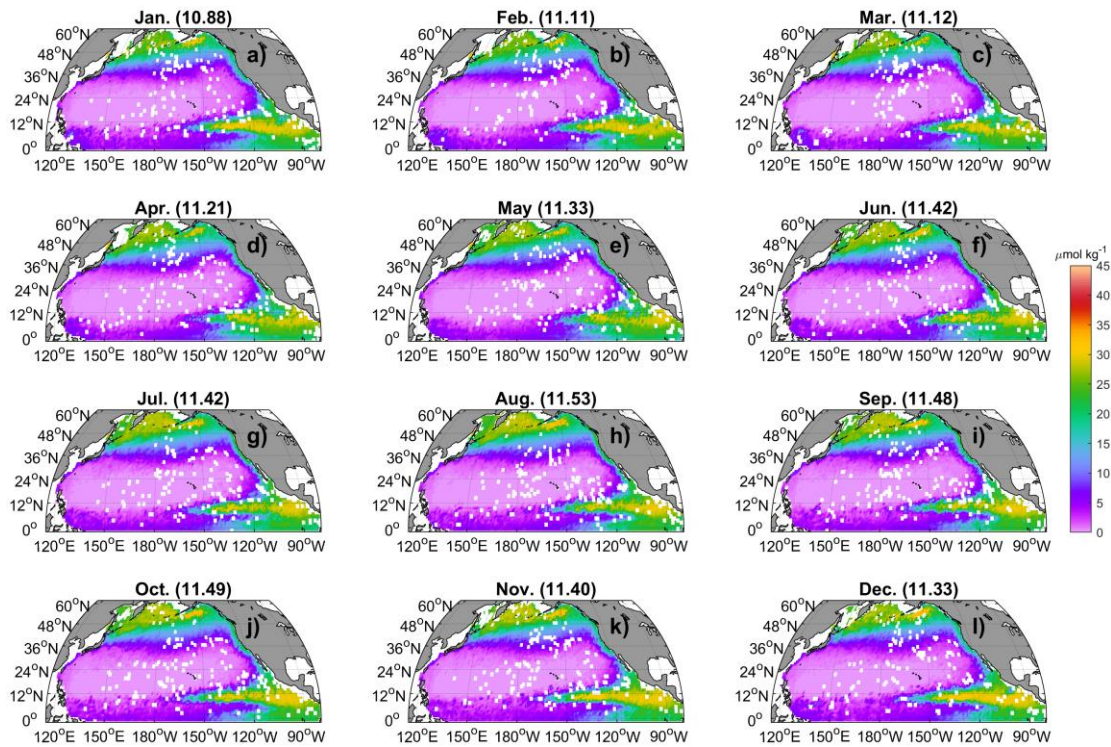
545 **Figure 10.** The monthly climatology of NO_x^- at 5 m in the North Pacific. Data are
 546 binned and averaged within $1 \times 1^\circ$ grid cells. The values in the title represent the spatial
 547 mean values.

548

549 At 100 m, NO_x^- concentrations are elevated particularly in the subarctic gyre, north
 550 of the Equator, and the eastern North Pacific, while the central regions, particularly the
 551 NPSG, exhibit lower values. At 500 m, NO_x^- concentrations display patterns similar to
 552 those at 100 m, except that the NO_x^- concentrations in the western NPSG are evidently
 553 lower than those in other regions (Fig. 13). At 1000 m, concentrations in the
 554 southwestern North Pacific Ocean are markedly lower than those in other regions (Fig.
 555 12). Below 100 m depth, seasonal variability in NO_x^- is minimal (Figs. 11–13).

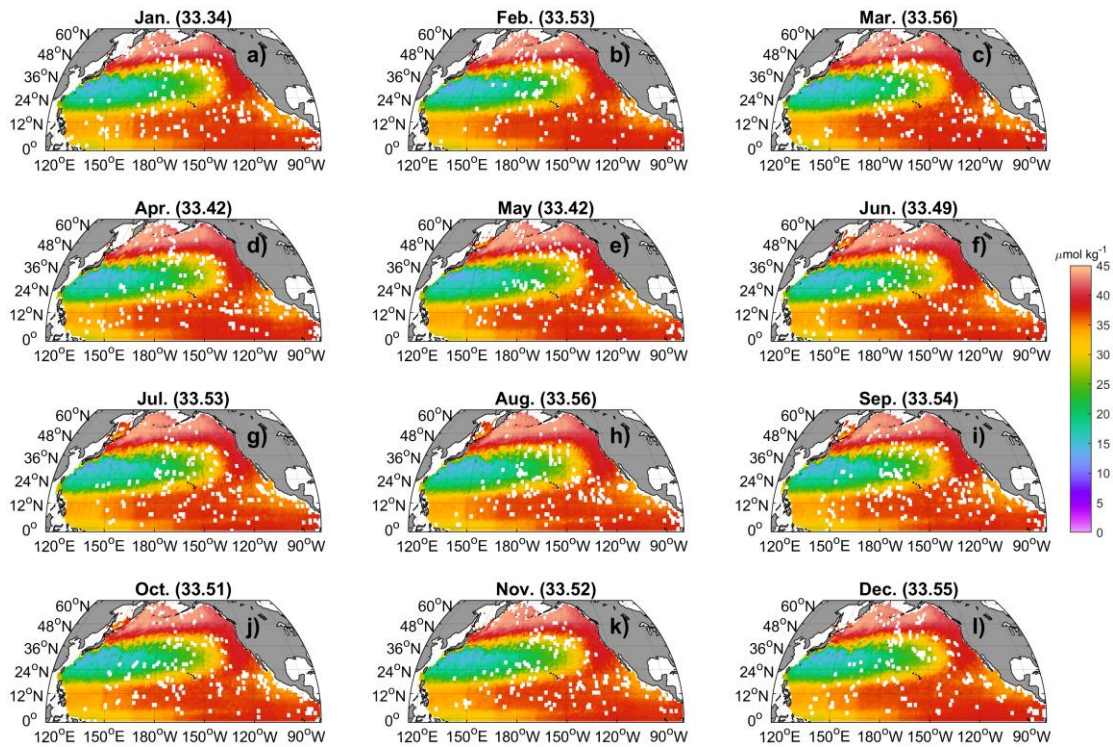
556 Compared to the World Ocean Atlas (WOA23) climatology (Figs. S17–S25), These
 557 results display patterns similar to WOA23 (Figs. S36–S44). The differences between
 558 the averaged values of these two climatologies are generally $<0.7 \mu\text{mol kg}^{-1}$ at the
 559 surface and $<1.5 \mu\text{mol kg}^{-1}$ at 100 m and 500 m. The maximum differences are found
 560 in July at a depth of 500 m (Figs. 13g and S38g). In that month and layer, WOA23
 561 shows a notably low mean NO_3^- value ($31.94 \mu\text{mol kg}^{-1}$) compared to its values in other

562 months (33.15 to 34.64 $\mu\text{mol kg}^{-1}$; Fig. S38) and compared to our climatology (33.34
 563 to 33.56 $\mu\text{mol kg}^{-1}$; Fig. 13). This discrepancy arises because the WOA23 climatology
 564 for July features a pronounced low- NO_3^- patch (down to 20 $\mu\text{mol kg}^{-1}$) within the
 565 eastern subarctic gyre, surrounded by waters with concentrations of $>35 \mu\text{mol kg}^{-1}$ (Fig.
 566 S38g). These regional differences are clearly visible in the difference maps between the
 567 two products (Figs. S45–S47). Generally although the seasonal patterns are similar in
 568 the surface layer, the reconstructed NO_x^- concentrations are lower than those in WOA23.
 569 In addition, our reconstructions capture finer spatial detail, exhibit less oversmoothing,
 570 and avoid artificial “bull’s-eye” patterns. It should be noted that our climatology is
 571 derived from the mean of existing data, which heavily relies on the spatiotemporal
 572 distribution of those data and may not represent the true climatological mean.



573 **Figure 11.** The monthly climatology of NO_x^- at 100 m in the North Pacific. Data are
 574 binned and averaged within $1 \times 1^\circ$ grid cells. The values in the title represent the spatial
 575 mean values.
 576

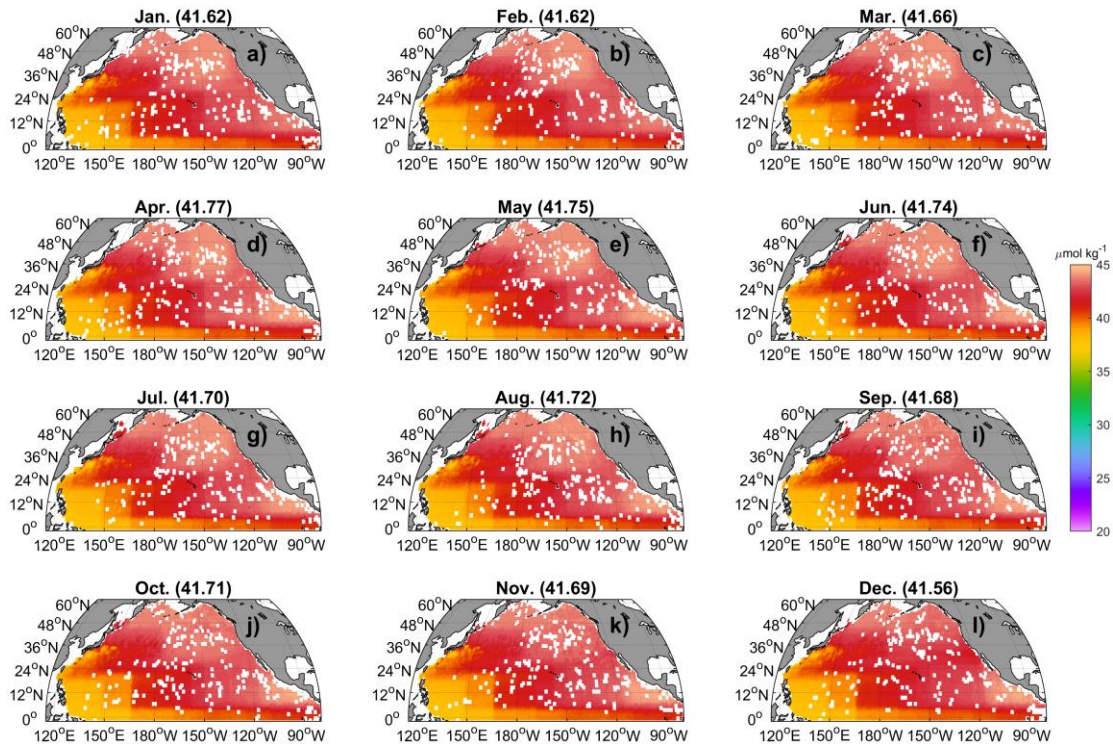
577



578

579 **Figure 12.** The monthly climatology of NO_x^- at 500 m in the North Pacific. Data are
 580 binned and averaged within $1 \times 1^\circ$ grid cells. The values in the title represent the spatial
 581 mean values.

582

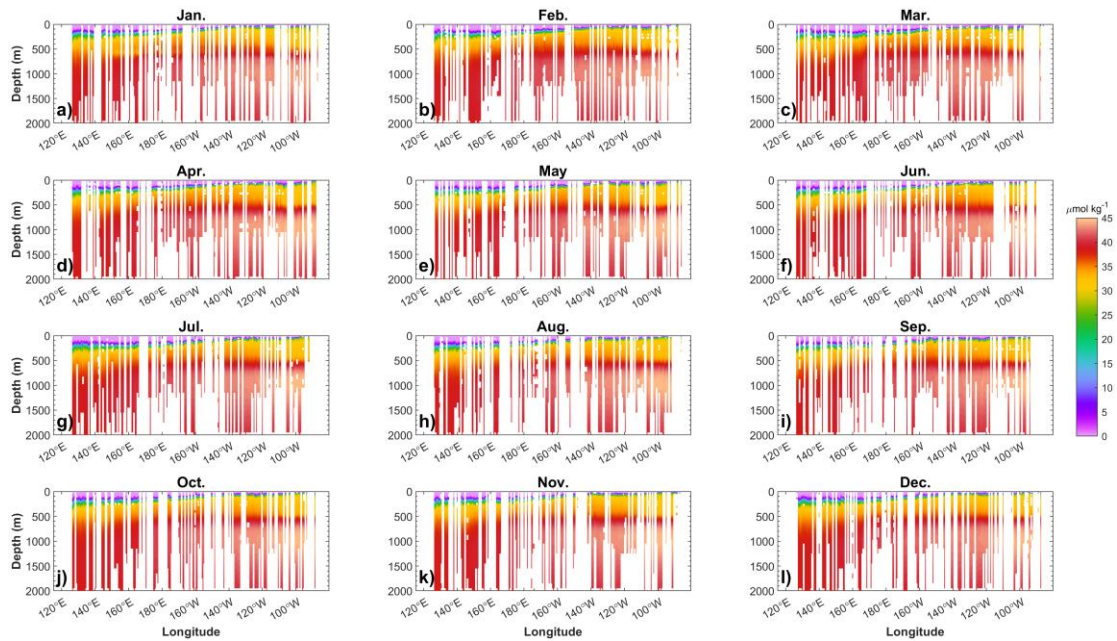


583

584 **Figure 13.** The monthly climatology of NO_x^- at 1000 m in the North Pacific. Data are
 585 binned and averaged within $1^\circ \times 1^\circ$ grid cells. The values in the title represent the spatial
 586 mean values.

587

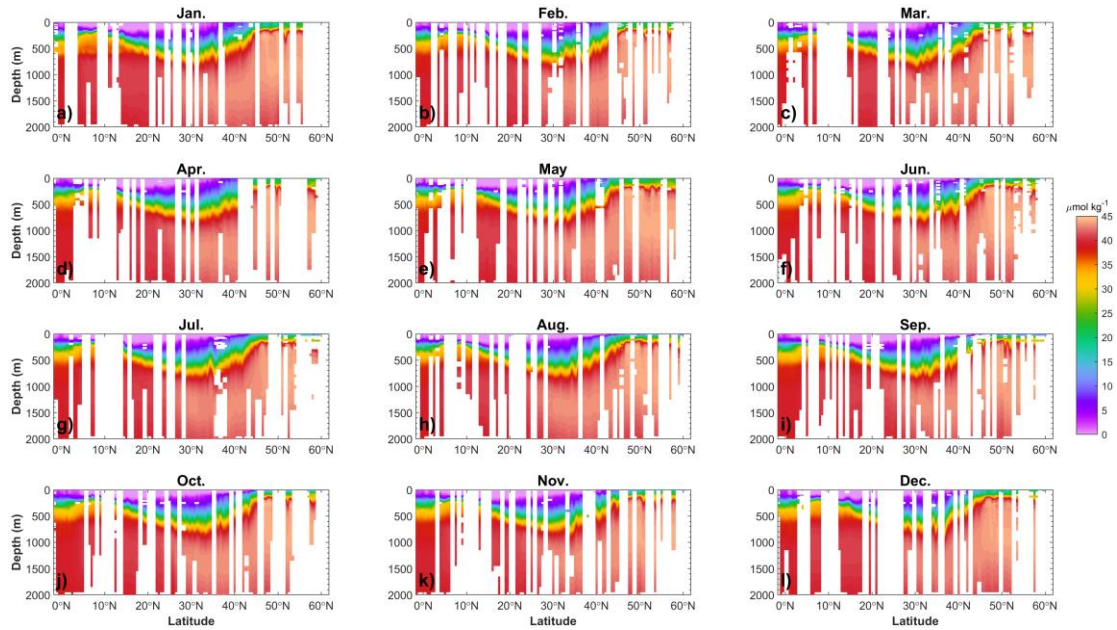
588 Sectional distributions of NO_x^- in the upper 2000 m along 10° N and 180° E were
 589 used as examples to illustrate the vertical profile distributions of nutrients within the
 590 North Pacific. At 10° N, NO_x^- concentrations increase from $\sim 0.0 \mu\text{mol kg}^{-1}$ at the
 591 surface to $\sim 45.0 \mu\text{mol kg}^{-1}$ at ~ 1000 m, followed by a decrease to $\sim 38.0 \mu\text{mol kg}^{-1}$ at
 592 2000 m. NO_x^- concentrations increase from west to the east in the North Pacific in the
 593 upper 300 m (Fig. 14). At 180° E, in the upper 500 m, meridional NO_x^- concentrations
 594 increase from the equator to the North Equatorial Current ($\sim 10^\circ$ N), decline within the
 595 subtropical gyre, and then increase toward the subarctic region (Fig. 15). Generally,
 596 seasonal differences of NO_x^- concentrations along both sections are not evident.



597

598 **Figure 14.** Zonal and monthly climatology of NO_x^- in the upper 2000 m at 10° N in the
 599 North Pacific. Data were binned and averaged within $1^\circ \times 1^\circ$ grid cells.

600



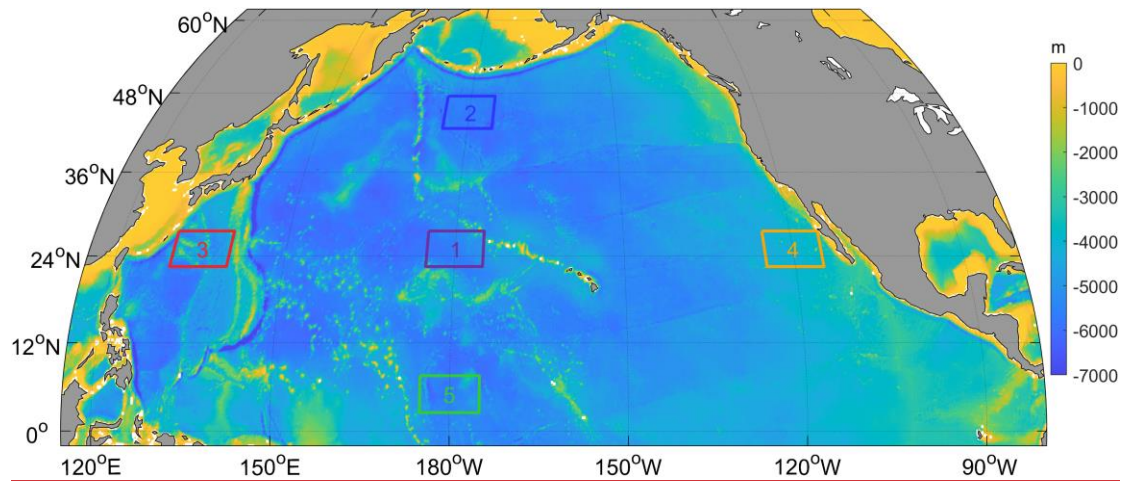
601

602 **Figure 15.** The monthly climatology of NO_x^- in the upper 2000 m at 170 °E section in
 603 the North Pacific. Data were binned and averaged within $1^\circ \times 1^\circ$ grid cells.

604

605 **3.4 Long-term variations of nutrients**

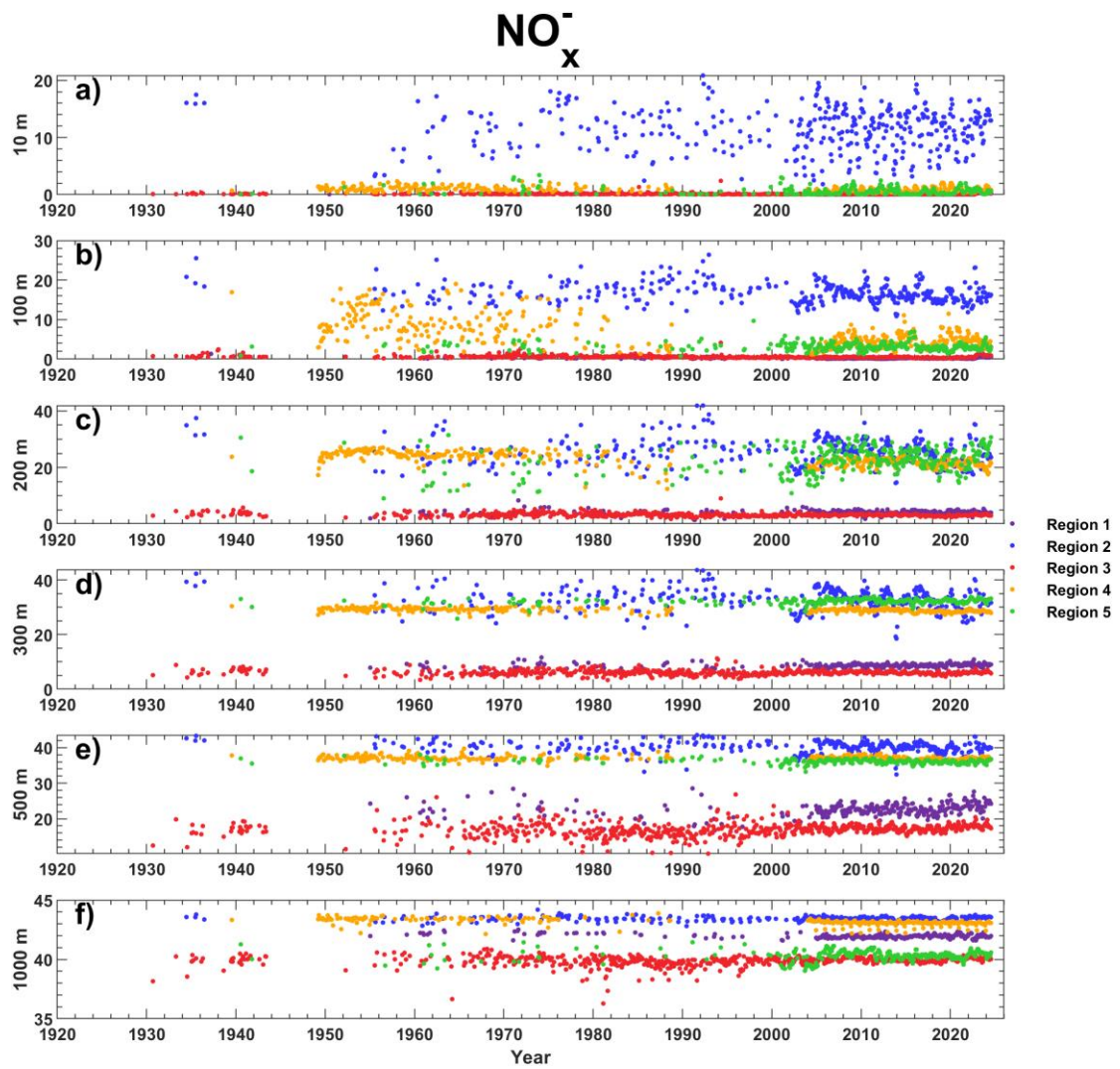
606 We present an initial analysis of long-term nutrient changes by examining five
 607 representative regions in the North Pacific, covering the subarctic gyre, the subtropical
 608 gyre, and equatorial areas (Fig. 16). The data are binned by region, month, and depth
 609 (10 m, 100 m, 200 m, 300 m, 500 m, and 1000 m) for regions 1–5. As shown in Fig. 17,
 610 these time series reveal notable interannual fluctuations of NO_3^- (with 2–5-year
 611 oscillations), providing a first-order view of low-frequency variability captured by the
 612 reconstruction. However, no evident long-term trend is found for nutrients. DIP and
 613 Si(OH)_4 display patterns similar to NO_3^- (Figs. S48–S49). In contrast, at depths of 200
 614 m and 300 m, NO_2^- displays an increasing trend in the central NPSG and a decreasing
 615 trend in the eastern NPSG during the 1970–2005 period (Fig. S50). More sophisticated
 616 trend analyses and basin-scale integrations are promising avenues for future work based
 617 on this newly reconstructed dataset.



618

619 **Figure 16.** Locations of five representative regions for analyzing long-term nutrient
 620 variations.

621



622

623 Figure 17. Time series of reconstructed NO₃⁻ concentrations at 10 m (a), 100 m (b),
624 200 m (c), 300 m (d), 500 m (e), and 1000 m (f) for regions 1–5 (see Fig. 16). Data
625 were binned by depth and region and then averaged by month.

627 **4 Data availability**

628 The database is available in a data repository (Du et al., 2025;
629 <https://zenodo.org/records/17140658>). Although the reconstruction results from RF,
630 LightGBM, and GPR are generally consistent, RF yields the best performance. To avoid
631 redundancy and minimize storage requirements—given the large volume of the data
632 files—only the nutrient data reconstructed by RF have been uploaded. Researchers may
633 contact the corresponding authors to request the reconstructions generated by
634 LightGBM and GPR.

636 **5 Conclusion**

637 In this study, we applied rigorous quality control procedures to clean hydrographic
638 and nutrient observations from CCHDO and WOD datasets. The cleaned CCHDO data
639 were then used to train three machine-learning models to relate nutrient concentrations
640 to spatial, temporal, and water-mass predictors. The models were applied to reconstruct
641 nutrient concentrations from hydrographic observations collected from WOD, though
642 most of which lack direct nutrient measurements. We assessed the model performance
643 using four data-partition strategies, and found that all models reproduced held-out data
644 with low RMSE_s-values. RF and GPR slightly outperformed LightGBM. The
645 application of these models to WOD hydrography yielded 472,652,680 reconstructed
646 nutrient concentrations across 1,920,634 stations and 35,744 cruises, spanning from
647 1895 to 2024. This represents a 2,127– to 2,393-fold increase compared to the original
648 volume of CCHDO nutrient data. The reconstruction captured the spatial, seasonal, and
649 interannual variations of water column nutrients in the North Pacific Ocean well.
650 Compared to the WOA23 climatology, the reconstruction-based nutrient climatology
651 exhibited more realistic spatial structures than WOA23. This high-quality and high-

652 resolution nutrient dataset ~~adds enables~~ historical nutrient estimation for locations and
653 times with ~~only solely~~ hydrographic measurements. The additionally potential
654 application of this dataset include: 1) investigating nutrient transport and budget in the
655 north Pacific; 2) spinning up and validating ocean biogeochemical models; 3) assessing
656 long-term nutrient trends driven by anthropogenic forcing and climate change; 4)
657 investigating nutrient stoichiometric changes and their ecological impacts under
658 climate variability. It also supports studies of climatological and long-term nutrient
659 variability under climate change and anthropogenic impacts, and provides transient
660 boundary conditions for ocean biogeochemical models in the Pacific Ocean.
661 Collectively, this resource facilitates advanced studies on marine biogeochemical
662 cycles, ecosystem dynamics, and climate-nutrient interactions.

663

664 **Author contributions**

665 CD and XL designed the study and dataset. CD, SK, MD, ZC, DS, and XL conceived
666 the project and secured the funding. CD, NZ, QL, ~~and~~HW and XL collected and
667 processed the data, developed the code, and performed the analysis. SK, MD, ZC, and
668 DS provided methodological guidance and advice. CD and NZ wrote the original draft.
669 All authors reviewed, edited the manuscript.

670

671 **Competing interests**

672 The contact author has declared that none of the authors has any competing interests.

673

674 **Acknowledgements**

675 This study was funded by the National Natural Science Foundation of China (Grants
676 42494885), National Key Research and Development Program of China
677 (Grant 2023YFF0805001), This study was funded by the National Natural Science
678 Foundation of China (Grants 42576215, 42494881), Innovational Fund for Scientific
679 and Technological Personnel of Hainan Province (Grant KJRC2023B04), and Natural
680 Science Foundation of Hainan Province (Grant 624MS037). We thank the CCHDO

681 (<https://cchdo.ucsd.edu/>) and the WOD ([https://www.ncei.noaa.gov/products/world-](https://www.ncei.noaa.gov/products/world-ocean-database)
682 [ocean-database](https://www.ncei.noaa.gov/products/world-ocean-database)) for providing the data used in this study. Special thanks are owed to all
683 scientists involved in data collection, analysis, and management for these programs.

684

685 **Declaration of generative AI and AI-assisted technologies in the writing process:**

686 During the preparation of this work the authors used deepseek to check the spelling and
687 grammar. After using this tool, the authors reviewed and edited the content as needed
688 and take full responsibility for the content of the publication.

689

690 **References**

691 Arteaga, L., Pahlow, M., and Oschlies, A.: Global monthly sea surface nitrate fields
692 estimated from remotely sensed sea surface temperature, chlorophyll, and
693 modeled mixed layer depth, *Geophys. Res. Lett.*, 42, 1130–1138, 2015.

694 Ascani, F., Richards, K. J., Firing, E., Grant, S., Johnson, K. S., Jia, Y., et al.: Physical
695 and biological controls of nitrate concentrations in the upper subtropical North
696 Pacific Ocean, *Deep-Sea Res. Pt. II*, 93, 119–134, 2013.

697 Barone, B., Church, M. J., Dugenne, M., Hawco, N. J., Jahn, O., White, A. E., et al.:
698 Biogeochemical dynamics in adjacent mesoscale eddies of opposite polarity,
699 *Global Biogeochem. Cy.*, 36, e2021GB007115, 2022.

700 Benitez-Nelson, C. R., Bidigare, R. R., Dickey, T. D., Landry, M. R., Leonard, C. L., et
701 al.: Mesoscale Eddies Drive Increased Silica Export in the Subtropical Pacific
702 Ocean, *Science*, 316, 1017–1021, 2007.

703 Bidigare, R. R., Chai, F., Landry, M. R., Lukas, R., Hannides, C. C. S., Christensen, S.
704 J., Karl, D. M., Shi, L., and Chao, Y.: Subtropical ocean ecosystem structure
705 changes forced by North Pacific climate variations, *J. Plankton Res.*, 31, 1131–
706 1139, 2009.

707 Bonnet, S., Caffin, M., Berthelot, H., and Moutin, T.: Hot spot of N₂ fixation in the
708 western tropical South Pacific pleads for a spatial decoupling between N₂ fixation
709 and denitrification, Proc. Natl. Acad. Sci. USA, 114, E2800–E2801, 2017.

710 Browning, T. J. and Moore, C. M.: Global analysis of ocean phytoplankton nutrient
711 limitation reveals high prevalence of co-limitation, Nat. Commun., 14, 5014, 2023.

712 Chelton, D. B., Schlax, M. G., Samelson, R. M., and de Szoeke, R. A.: Global
713 observations of large oceanic eddies, Geophys. Res. Lett., 34, L15606, 2007.

714 Chen, S., Hu, C., Barnes, B. B., Wanninkhof, R., Cai, W., Barbero, L., and Pierrot, D.:
715 A machine learning approach to estimate surface ocean *p*CO₂ from satellite
716 measurements, Remote Sens. Environ., 228, 203–226, 2019.

717 Chen, S., Meng, Y., Lin, S., Yu, Y., and Xi, J.: Estimation of sea surface nitrate from
718 space: Current status and future potential, Sci. Total Environ., 899, 165690, 2023.

719 Chen, S., Meng, Y., Shang, S., Zheng, M., Wang, Y., and Chai, F.: Remote estimates of
720 sea surface nitrate and its trends from ocean color in the northwest Pacific, J.
721 Geophys. Res., 129, e2023JC019846, 2024.

722 Dai, M., Luo, Y., Achterberg, E. P., Browning, T. J., Cai, Y., Cao, Z., Chai, F., Chen, B.,
723 Church, M. J., Ci, D., Du, C., Gao, K., Guo, X., Hu, Z., Kao, S., Laws, E. A., Lee,
724 Z., Lin, H., Liu, Q., et al.: Upper Ocean biogeochemistry of the oligotrophic North
725 Pacific subtropical gyre: From nutrient sources to carbon export, Rev. Geophys.,
726 61, e2022RG000800, 2023.

727 Du, C., Zheng, N., Kao, S.-J., Dai, M., Cao, Z., Shi, D., Li, Q., Wang, H., and Li, X.:
728 Validated temperature and salinity data, and reconstructed nutrient concentrations
729 in the North Pacific (1895 – 2024). Zenodo, <https://zenodo.org/records/17451417>,
730 2025.

731 Dave, A. C. and Lozier, M. S.: Local stratification control of marine productivity in the
732 subtropical North Pacific, J. Geophys. Res., 115, C12032, 2010.

733 Deutsch, C. and Weber, T.: Nutrient Ratios as a Tracer and Driver of Ocean
734 Biogeochemistry, Annu. Rev. Mar. Sci., 4, 113–138, 2012.

735 Dong, L., Qi, J., Yin, B., Zhi, H., Li, D., Yang, S., Wang, W., Cai, H., and Xie, B.:
736 Reconstruction of subsurface salinity structure in the South China Sea using
737 satellite observations: a LightGBM-Based Deep forest method, *Remote Sens.*, 14,
738 3494, 2022.

739 Du, C., He, R., Liu, Z., Huang, T., Wang, L., Yuan, Z., Xu, Y., Wang, Z., and Dai, M.:
740 Climatology of nutrient distributions in the South China Sea based on a large data
741 set derived from a new algorithm, *Prog. Oceanogr.*, 195, 102586, 2021.

742 Dugdale, R. C., Morel, A., Bricaud, A., and Wilkerson, F. P.: Modeling new production
743 in upwelling centers: A case study of modeling new production from remotely-
744 sensed temperature and color, *J. Geophys. Res.*, 94, 18119–18132, 1989.

745 [Eugster, O. and Gruber, N.: A probabilistic estimate of global marine N-fixation and](#)
746 [denitrification, *Glob. Biogeochem. Cycles*, 26, GB4013, 2012.](#)

747 Fuhr, M., Laukert, G., Yu, Y., Nürnberg, D., and Frank, M.: Tracing water mass mixing
748 from the Equatorial to the North Pacific Ocean with dissolved neodymium
749 isotopes and concentrations, *Front. Mar. Sci.*, 7, 603761, 2021.

750 Gan, M., Pan, S., Chen, Y., Cheng, C., Pan, H., and Zhu, X.: Application of the Machine
751 Learning LightGBM model to the prediction of the water levels of the Lower
752 Columbia River, *J. Mar. Sci. Eng.*, 9, 496, 2021.

753 Garcia, H. E., Boyer, T. P., Locarnini, R. A., Reagan, J. R., Mishonov, A. V., Baranova,
754 O. K., Paver, C. R., Wang, Z., Bouchard, C. N., Cross, S. L., Seidov, D., and
755 Dukhovskoy, D.: *World Ocean Database 2023: User’s Manual*. A.V. Mishonov,
756 Technical Ed., NOAA Atlas NESDIS, 98, 129 pp., 2024.

757 Goes, J. I., Saino, T., Oaku, H., and Jiang, D. L.: A Method for Estimating Sea Surface
758 Nitrate Concentrations from Remotely Sensed SST and Chlorophyll - A Case
759 Study for the North Pacific Ocean Using OCTS/ADEOS Data, *IEEE Trans. Geosci.*
760 *Remote Sens.*, 37, 1633–1644, 1999.

761 Hu, C., Feng, L., and Guan, Q.: A machine learning approach to estimate surface
762 chlorophyll *a* concentrations in global oceans from satellite measurements, *IEEE*
763 *Trans. Geosci. Remote Sens.*, 59, 4590–4607, 2021.

764 Huang, Y., Nicholson, D., Huang, B., and Cassar, N.: Global estimates of marine gross
765 primary production based on machine learning upscaling of field observations,
766 *Global Biogeochem. Cy.*, 35, e2020GB006718, 2021.

767 Huang, Y., Tagliabue, A., and Cassar, N.: Data-Driven Modeling of Dissolved Iron in
768 the Global Ocean, *Front. Mar. Sci.*, 9, 837183, 2022.

769 Kamykowski, D., Zentara, S.-J., Morrison, J. M., and Switzer, A. C.: Dynamic global
770 patterns of nitrate, phosphate, silicate, and iron availability and phytoplankton
771 community composition from remote sensing data, *Global Biogeochem. Cy.*, 16,
772 1077, 2002.

773 Kamykowski, D.: A preliminary model of the relationship between temperature and
774 plant nutrients in the upper ocean, *Deep-Sea Res.*, 34, 1067–1079, 1987.

775 Kamykowski, D.: Estimating upper ocean phosphate concentrations using ARGO float
776 temperature profiles, *Deep-Sea Res. Pt. I*, 55, 1580–1589, 2008.

777 Karl, D. M. and Church, M. J.: Ecosystem structure and dynamics in the North Pacific
778 Subtropical Gyre: new views of an old ocean, *Ecosystems*, 20, 433–457, 2017.

779 Karl, D. M., Letelier, R. M., Bidigare, R. R., Björkman, K. M., Church, M. J., Dore, J.
780 E., and White, A. E.: Seasonal-to-decadal scale variability in primary production
781 and particulate matter export at Station ALOHA, *Prog. Oceanogr.*, 195, 102563,
782 2021.

783 Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., and Liu, T.-Y.:
784 Lightgbm: A highly efficient gradient boosting decision tree, *Adv. Neural Inf.*
785 *Process. Syst.*, 30, 3147–3155, 2017.

786 Lee, G. S., Lee, J. H., and Cho, H. Y.: Spatiotemporal estimation of nutrient data from
787 the northwest pacific and east Asian seas, *Sci. Data*, 10, 2023.

788 Liaw, A. and Wiener, M.: Classification and regression by randomForest, *R News*, 2,
789 18–22, 2002.

790 Lipschultz, F., Bates, N. R., Carlson, C. A., and Hansell, D. A.: New production in the
791 Sargasso Sea: History and current status, *Global Biogeochem. Cy.*, 16, 1001, 2002.

792 Liu, H., Lin, L., Wang, Y., Du, L., Wang, S., Zhou, P., Yu, Y., Gong, X., and Lu, X.:
793 Reconstruction of Monthly Surface Nutrient Concentrations in the Yellow and
794 Bohai Seas from 2003–2019 Using Machine Learning, *Remote Sens.*, 14, 5021,
795 2022.

796 Madani, N., Parazoo, N. C., Manizza, M., Chatterjee, A., Carroll, D., Menemenlis, D.,
797 Fouest, V. L., Matsuoka, A., Luis, K. M., Serra-Pompei, C., and Miller, C. E.: A
798 machine learning approach to produce a continuous Solar-Induced chlorophyll
799 fluorescence over the Arctic Ocean, *J. Geophys. Res. Machine Learn. Comput.*, 1,
800 2024.

801 [Martino, M., Hamilton, D. S., Baker, A. R., Jickells, T., Bromley, T., Nojiri, Y., Quack,](#)
802 [B., and Boyd, P. W.: Western Pacific atmospheric nutrient deposition fluxes, their](#)
803 [impact on surface ocean productivity, *Glob. Biogeochem. Cycles*, 28, 712–728,](#)
804 [2014.](#)

805 Mishonov, A. V., Boyer, T. P., Baranova, O. K., Bouchard, C. N., Cross, S. L., Garcia,
806 H. E., Locarnini, R. A., Paver, C. R., Reagan, J. R., Wang, Z., Seidov, D., Grodsky,
807 A. I., and Beauchamp, J. G.: *World Ocean Database 2023*, C. Bouchard, Technical
808 Ed., NOAA Atlas NESDIS, 97, 2024.

809 Moore, C. M., Mills, M. M., Arrigo, K. R., Berman - Frank, I., Bopp, L., Boyd, P. W.,
810 Galbraith, E. D., Geider, R. J., Guieu, C., Jaccard, S. L., Jickells, T. D., Lenton, T.
811 M., Mahowald, N. M., Marañón, E., Marinov, I., Moore, J. K., Nakatsuka, T.,
812 Oschlies, A., Saito, M. A., Thingstad, T., Tsuda, A., and Ulloa, O.: Processes and
813 patterns of oceanic nutrient limitation, *Nat. Geosci.*, 6, 701–710, 2013.

814 Możejko, J. and Gniot, R.: Application of Neural Networks for the Prediction of Total
815 Phosphorus Concentrations in Surface Waters, *Pol. J. Environ. Stud.*, 17, 363–368,
816 2008.

817 Palacios, D. M., Hazen, E. L., Schroeder, I. D., and Bograd, S. J.: Modeling the
818 temperature-nitrate relationship in the coastal upwelling domain of the California
819 Current, *J. Geophys. Res.*, 118, 1–17, 2013.

820 [Qi, J., Yu, Y., Yao, X., Yuan, G., and Gao, H.: Dry deposition fluxes of inorganic](#)
821 [nitrogen and phosphorus in atmospheric aerosols over the Marginal Seas and](#)
822 [Northwest Pacific, Atmos. Res., 245, 105076, 2020.](#)

823 Reagan, J. R., Boyer, T. P., García, H. E., Locarnini, R. A., Baranova, O. K., Bouchard,
824 C., Cross, S. L., Mishonov, A. V., Paver, C. R., Seidov, D., Wang, Z., and
825 Dukhovskoy, D.: World Ocean Atlas 2023, NOAA National Centers for
826 Environmental Information, Dataset, NCEI Accession 0270533, 2024.

827 Sarangi, P. K., Thangaradjou, T., Kumar, A. S., and Balasubramanian, T.: Development
828 of nitrate algorithm for the southwest bay of bengal water and its implication using
829 remote sensing satellite datasets, IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.,
830 4, 983–991, 2011.

831 Sigman, D. M. and Hain, M. P.: The Biological Productivity of the Ocean, Nat. Educ.
832 Knowl., 3, 21, 2012.

833 Steinhoff, T., Friedrich, T., Hartman, S. E., Oschlies, A., Wallace, D. W. R., and
834 Körtzinger, A.: Estimating mixed layer nitrate in the North Atlantic Ocean,
835 Biogeosciences, 7, 795–807, 2010.

836 Su, H., Lu, X., Chen, Z., Zhang, H., Lu, W., and Wu, W.: Estimating Coastal
837 Chlorophyll-A Concentration from Time-Series OLCI Data Based on Machine
838 Learning, Remote Sens., 13, 576, 2021.

839 Sundararaman, H. K. K. and Shanmugam, P.: Estimates of the global ocean surface
840 dissolved oxygen and macronutrients from satellite data, Remote Sens. Environ.,
841 311, 114243, 2024.

842 Switzer, A. C., Kamykowski, D., and Zentara, S.-J.: Mapping nitrate in the global ocean
843 using remotely sensed sea surface temperature, J. Geophys. Res., 108, 345–359,
844 2003.

845 Talley, L. D., Pickard, G. L., Emery, W. J., and Swift, J. H.: Descriptive Physical
846 Oceanography, An Introduction, Sixth Edition, Academic Press, 350–362 pp.,
847 2011.

- 848 Wang, C., Su, B., Sun, J., Hu, X., and Liu, J.: A regional ocean database for the Coastal
849 China Sea. *Sci Data*, 12, 1550, 2025.
- 850 Wang, L., Xu, Z., Gong, X., Zhang, P., Hao, Z., You, J., Zhao, X., and Guo, X.:
851 Estimation of nitrate concentration and its distribution in the northwestern Pacific
852 Ocean by a deep neural network model, *Deep Sea Res. I*, 195, 104005, 2023.
- 853 Wang, Z., Wang, G., Guo, X., Hu, J., and Dai, M.: Reconstruction of High-Resolution
854 Sea Surface Salinity over 2003–2020 in the South China Sea Using the Machine
855 Learning Algorithm LightGBM Model, *Remote Sens.*, 14, 6147, 2022.
- 856 Yang, G. G., Wang, Q., Feng, J., He, L., Li, R., Lu, W., Liao, E., and Lai, Z.: Can three-
857 dimensional nitrate structure be reconstructed from surface information with
858 artificial intelligence? – A proof-of-concept study, *Sci. Total Environ.*, 924,
859 171365, 2024.
- 860 Yasunaka, S., Mitsudera, H., Whitney, F., and Nakaoka, S.: Nutrient and dissolved
861 inorganic carbon variability in the North Pacific, *J. Oceanogr.*, 77, 3–16, 2021.
- 862 Yasunaka, S., Nojiri, Y., Nakaoka, S., Ono, T., Whitney, F. A., and Telszewski, M.:
863 Mapping of sea surface nutrients in the North Pacific: Basin-wide distribution and
864 seasonal to interannual variability, *J. Geophys. Res. Oceans*, 119, 7756–7771,
865 2014.
- 866 Yasunaka, S., Ono, T., Nojiri, Y., Whitney, F. A., Wada, C., Murata, A., Nakaoka, S.,
867 and Hosoda, S.: Long-term variability of surface nutrient concentrations in the
868 North Pacific, *Geophys. Res. Lett.*, 43, 3389–3397, 2016.
- 869 Yu, X. R., Wen, Z., Jiang, R., Yang, J.-Y. T., Cao, Z., Hong, H., Zhou, Y., and Shi, D.:
870 Assessing N₂ fixation flux and its controlling factors in the (sub)tropical western
871 North Pacific through high-resolution observations, *Limnol. Oceanogr. Lett.*, 9,
872 716 – 724, 2024.
- 873 Zhong, A., Wang, D., Gong, F., Zhu, W., Fu, D., Zheng, Z., Huang, J., He, X., and Bai,
874 Y.: Remote sensing estimates of global sea surface nitrate: Methodology and
875 validation, *Sci. Total Environ.*, 950, 175362, 2024.

876

