

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 large database of ~473 million nutrient data points across 1.92
11 million stations (1895–2024), expanding the number of nutrient data points by a
12 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 $<3.07 \mu\text{mol}$
31 kg^{-1} for NO_3^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$, respectively. The validated models were then
32 applied to reconstruct nutrient concentrations from the hydrographic observations in
33 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. The
40 dataset generated in this study is openly available on Zenodo at
41 <https://zenodo.org/records/17451417>.

42

43 **1 Introduction**

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

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

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

72 (2008) used Artificial Neural Networks (ANNs) to model time series of total
73 phosphorus concentrations in the Odra River. Self-organizing maps (SOMs) were used
74 to estimate mixed layer nitrate and sea surface nutrients in the open ocean (Steinhoff et
75 al., 2010; Yasunaka et al., 2014). Liu et al. (2022) applied Support Vector Regression,
76 Random Forest Regression, and ANNs to reconstruct monthly surface nutrient
77 concentrations in the Yellow and Bohai Seas from 2003 to 2019. Their results revealed
78 pronounced seasonal and spatial variability in nutrient levels and underscored the
79 influence of environmental drivers such as sea surface temperature and salinity.
80 Similarly, Sundararaman and Shanmugam (2024) employed Gaussian Process
81 Regression (GPR) models to estimate global ocean surface macronutrient
82 concentrations using satellite-derived data, achieving high accuracy and demonstrating
83 their suitability for large-scale marine ecosystem monitoring. Yang et al. (2024)
84 employed a U-net and Earthformer to reconstruct the three-dimensional nitrate
85 distribution by integrating surface data including wind speed, sea surface temperature,
86 chlorophyll *a*, solar radiation, and precipitation in the Indian Ocean. These
87 advancements highlight the expanding role of machine learning in marine biochemical
88 data fusion and provide novel insights into nutrient dynamics and their ecological
89 impacts.

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

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

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

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

129 chlorophyll *a*, and mixed layer depth. These maps revealed seasonal and interannual
130 variability in surface nutrient distributions in the northern North Pacific. To investigate
131 long-term changes, Yasunaka et al. (2016) applied Optimal Interpolation to analyze the
132 spatial and temporal evolution of surface nutrient concentrations. Lee et al. (2023)
133 provided spatiotemporally gridded nitrate and phosphate data in the northwest Pacific
134 from 1980 to 2019 using the spatiotemporal kriging technique. Wang et al. (2023) used
135 the deep neural network model to estimate nitrate concentrations in the upper
136 northwestern Pacific Ocean using temperature and salinity as the primary input
137 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 Table 1. Information on nutrients and their associated hydrographic data collected
 159 from CLIVAR and Carbon Hydrographic Data Office (CCHDO) and the information
 160 after quality control (QC).

	Original data information		Data information after QC	
	Data	Stations	Data	Stations
Temperature	327792	15127	327688	15125
Salinity	328502	15274	328275	15269
NO _x ⁻	217725	9588	213962	9021
NO ₂ ⁻	197617	8233	197539	8228
DIP	225403	9623	222234	9474
Si(OH) ₄	212660	8220	210447	8121

161 Hydrographic data for nutrient reconstruction were obtained from the World Ocean
 162 Database (WOD; Mishonov et al., 2024), which compiles observations from various
 163 platforms, including Autonomous Pinniped Bathythermograph (APB), Conductivity-
 164 Temperature-Depth profiler (CTD), Drifting Buoy (DRB), Glider (GLD), Mechanical
 165 Bathythermograph (MBT), Moored Buoy (MRB), Ocean Station Data (OSD), Profiling
 166 Float (PFL), and Undulating Oceanographic Recorder (UOR). Since nutrient
 167 reconstruction models rely on relationships with water masses, only samples containing
 168 both temperature and salinity measurements were used; therefore, most APB
 169 observations, which record only temperature, were excluded. Among these platforms,
 170 CTD, OSD, and PFL provided the majority of usable data. Additionally, several
 171 marginal seas—including the South China Sea, the Yellow Sea, the Sea of Japan, and
 172 the Sea of Okhotsk—were excluded from this study because they are semi-enclosed
 173 and strongly influenced by terrestrial inputs. The spatial domain was consistent with
 174 that used for the CCHDO dataset, while the temporal coverage extended from 1875 to
 175 2024. In total, 577,215,683 data points from 2,284,448 stations across 40,113 original
 176 cruises were collected (Table 2). In addition, the OSD data before 1970 were extracted
 177 for nutrient validation in Section 3.1. A total of 102,424, 125,142, 447,335, and 294,734
 178 data points were collected for NO₃⁻, NO₂⁻, DIP, and Si(OH)₄, respectively.

179 Table 2. Information on hydrographic data collected from World Ocean Database, and
 180 the data information after quality control (QC). See main text for acronyms' full
 181 names.

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

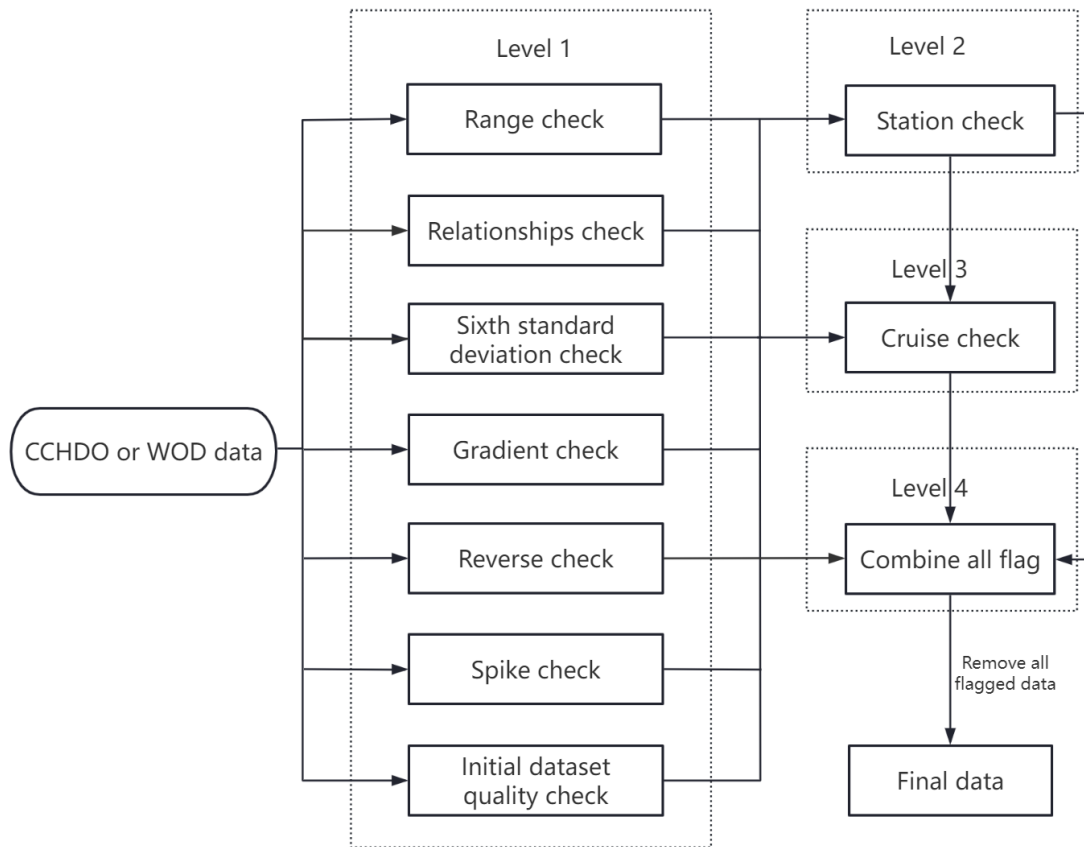
182

183 2.2 Data quality control

184 Given that the data were collected from multiple platforms using various methods
 185 over a long-time span and broad spatial range, quality control (QC) was essential (Du
 186 et al., 2021; Wang et al., 2025). Following the QC procedures developed by the World
 187 Ocean Database (WOD) (Garcia et al., 2024), we applied comprehensive QC protocols
 188 (Fig. 1) to both CCHDO and WOD datasets, including hydrographic and nutrient
 189 variables.

190 Four levels of QC were applied to identify and remove potentially erroneous or low-
 191 quality records from the CCHDO and WOD datasets. The first level targeted individual
 192 measurements, including several checks. (1) A range check was conducted by defining
 193 depth-dependent acceptable value ranges for each parameter; data falling outside these
 194 ranges were flagged as invalid. This check was applied to temperature, salinity, NO_x^- ,
 195 NO_2^- , DIP, and $\text{Si}(\text{OH})_4$. Note that the NO_x^- denotes the sum concentration of NO_2^- and

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



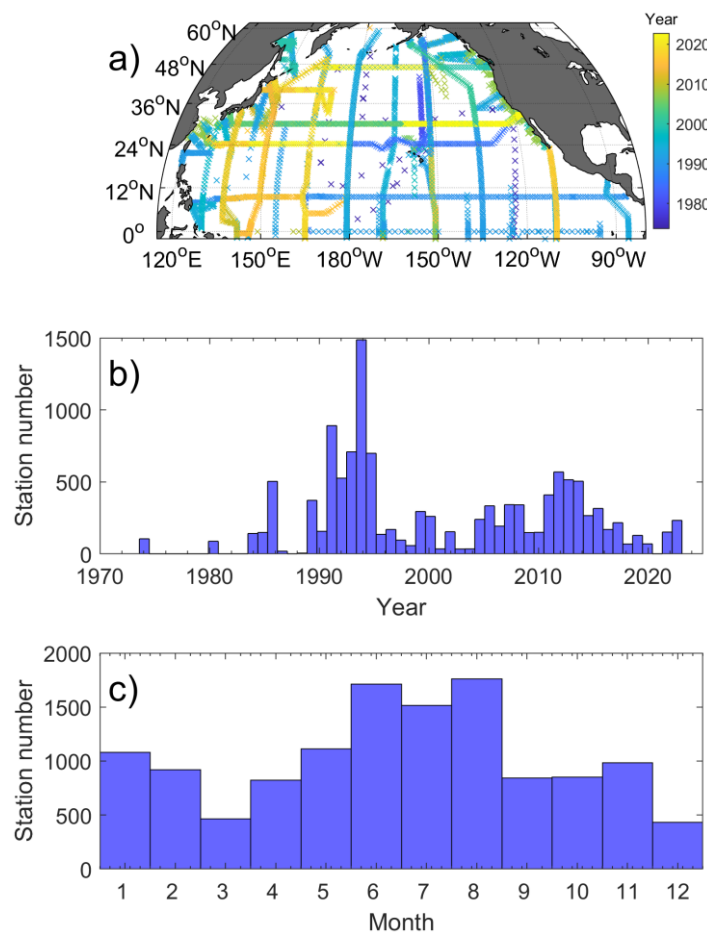
218

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

222

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

232 After quality control, the CCHDO dataset retained 214,943 (9,120), 197,539 (8,228),
 233 222,234 (9,457) and 210,447 (8,123) data points (stations), accounting for 94.2%
 234 (95.1%), 100.0% (99.9%), 98.6% (98.5%) and 99.0% (98.8%) of the original data
 235 points (stations) for NO_x^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$, respectively (Table 1). The retained
 236 stations cover nearly the entire North Pacific Ocean (Fig. 2a), spanning from 1972 to
 237 2023. Most observations were collected after 1980, with a substantial increase after
 238 1990 (Fig. 2b). Seasonally, the number of stations in June, July, and August was
 239 approximately three times greater than that in March and December (Fig. 2c).



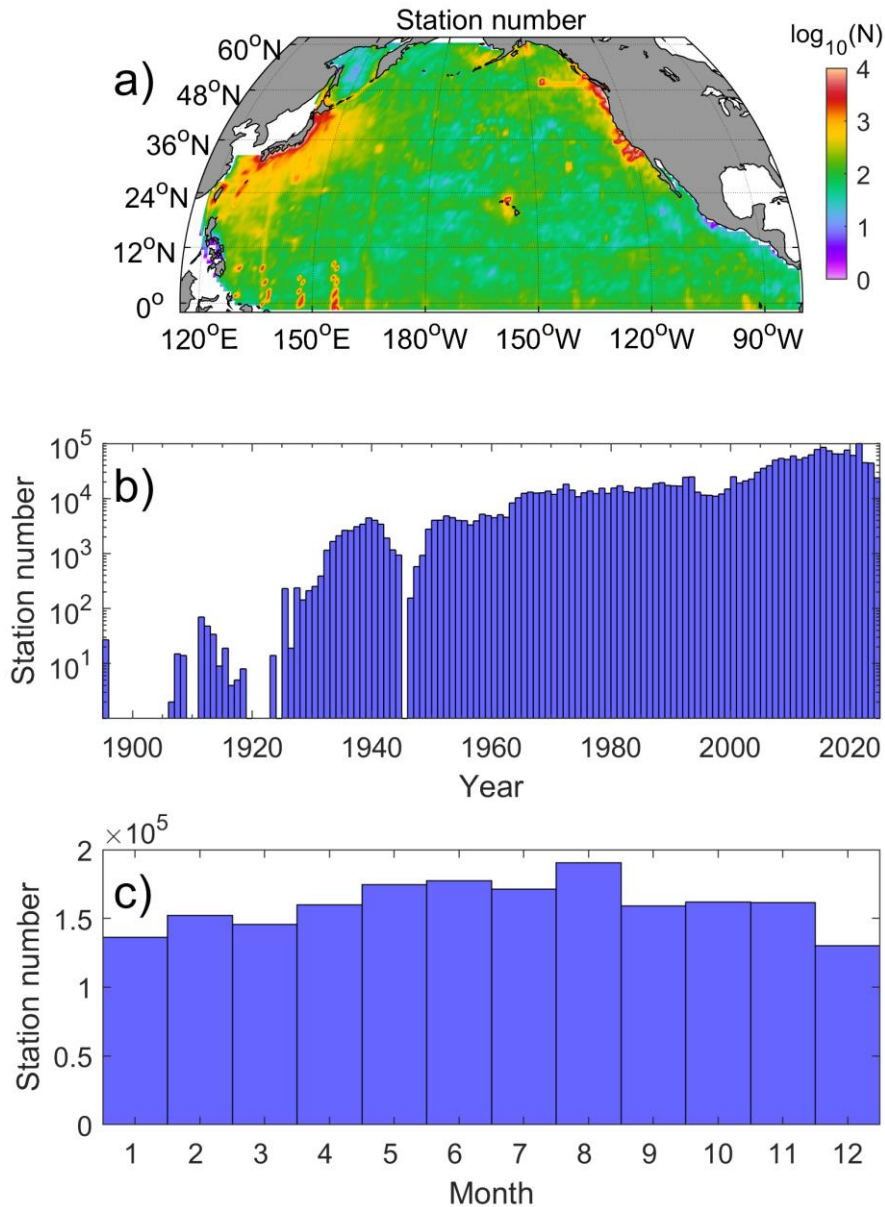
240

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

244

245 Following quality control, the final WOD dataset comprised 472,652,680
 246 temperature and salinity data points from 1,920,634 stations across 35,744 cruises,

247 spanning 1895 to 2024. These represent 81.9% of the original observations, 84.1% of
248 the original stations, and 89.1% of the original cruises, respectively (Table 2). Spatially,
249 station counts per $1^{\circ}\times 1^{\circ}$ grid cell range from 1 to 31,851, with a mean of 249 stations
250 per cell (Fig. 3a). High sampling densities are found off eastern Japan and western
251 North America, resulting from high frequency observations from CTD and OSD
252 platforms, whereas elevated counts in the southwestern North Pacific primarily result
253 from MRB observations. Temporally, fewer than 300 stations per year were collected
254 before 1930. The annual number of stations exceeded 10,000 after 1964 and peaked at
255 approximately 100,000 in 2021 (Fig. 3b). Seasonally, station numbers are highest from
256 May to August (Fig. 3c). Overall, the collected WOD dataset provides 2127–2393 times
257 more observations and 202 times more station records than the CCHDO dataset.



258

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

262

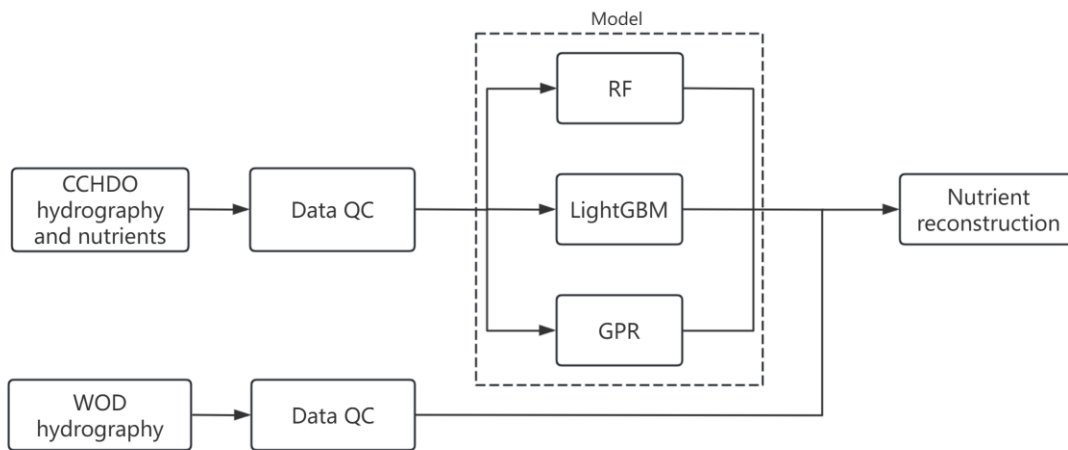
263 **2.3 Machine learning and nutrient reconstruction**

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

268 concentrations. These methods are widely used in marine science (Hu et al., 2021;
269 Huang et al., 2022; Yu et al., 2024; Chen et al., 2023; Sundararaman and Shanmugam,
270 2024). The use of diverse models helps reduce algorithm selection bias. RF is an
271 ensemble technique based on bagging, which builds multiple independent decision
272 trees and aggregates their outputs by voting or averaging (Liaw and Wiener, 2002). Its
273 strengths include high predictive accuracy and reduced overfitting owing to the large
274 number of trees. RF has been applied to predict global primary production (Huang et
275 al., 2021), chlorophyll concentrations (Madani et al., 2024), nutrients (Chen et al., 2023;
276 Chen et al., 2024), dissolved iron (Huang et al., 2022), surface ocean $p\text{CO}_2$ (Chen et al.,
277 2019), and N_2 fixation rates (Yu et al., 2024).

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

290

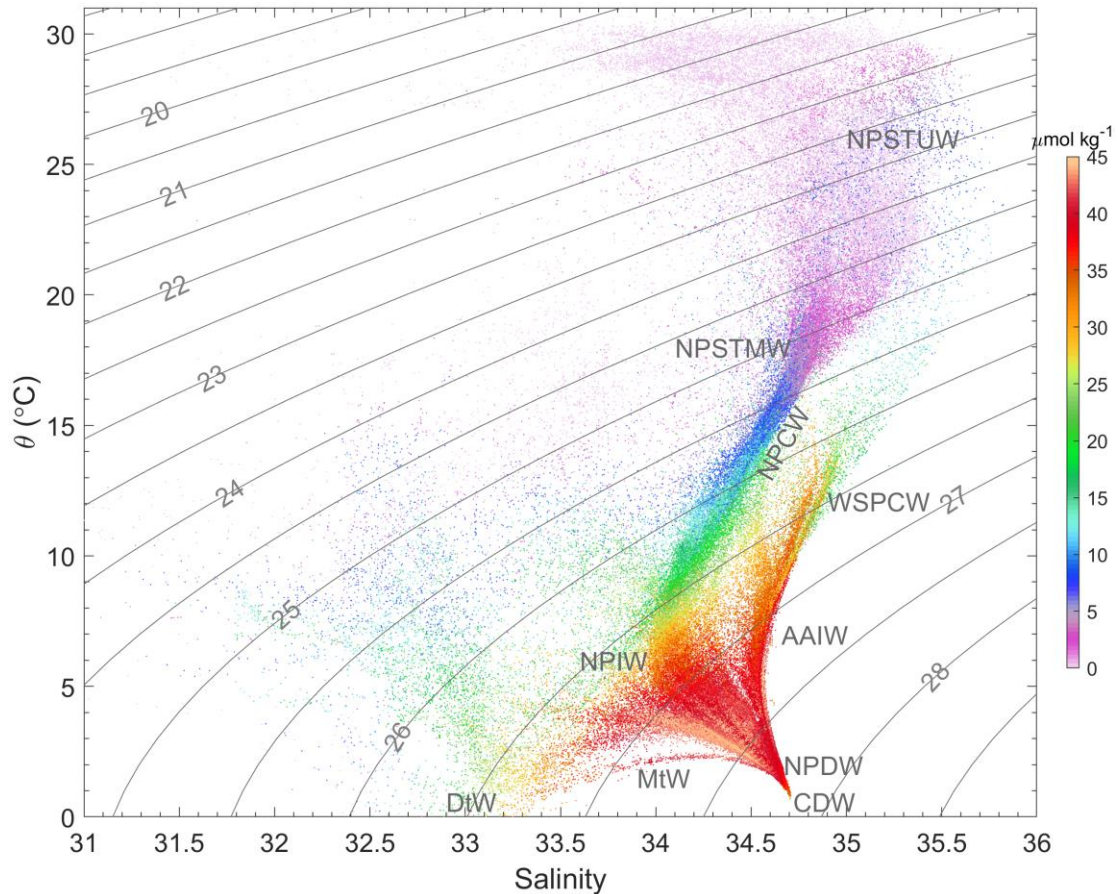


291

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

294

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



307

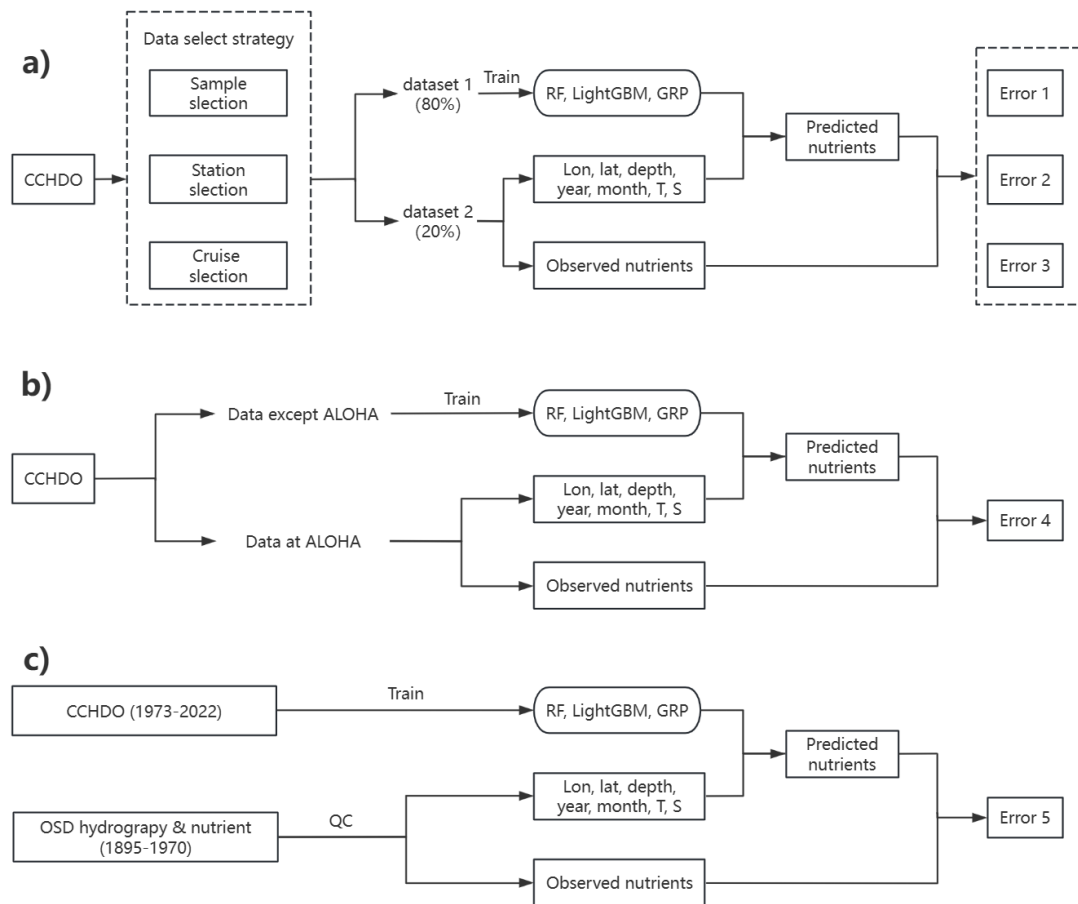
308 **Figure 5.** The water masses (indicated by salinity and potential temperature (θ)) and
 309 NO_x^- ($\text{NO}_3^- + \text{NO}_2^-$; color shading) relationships in the North Pacific. The temperature
 310 and salinity data were collected from the CCHDO dataset. The gray contour lines and
 311 number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

319

320 **3 Results**

321 **3.1 Error estimation**

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



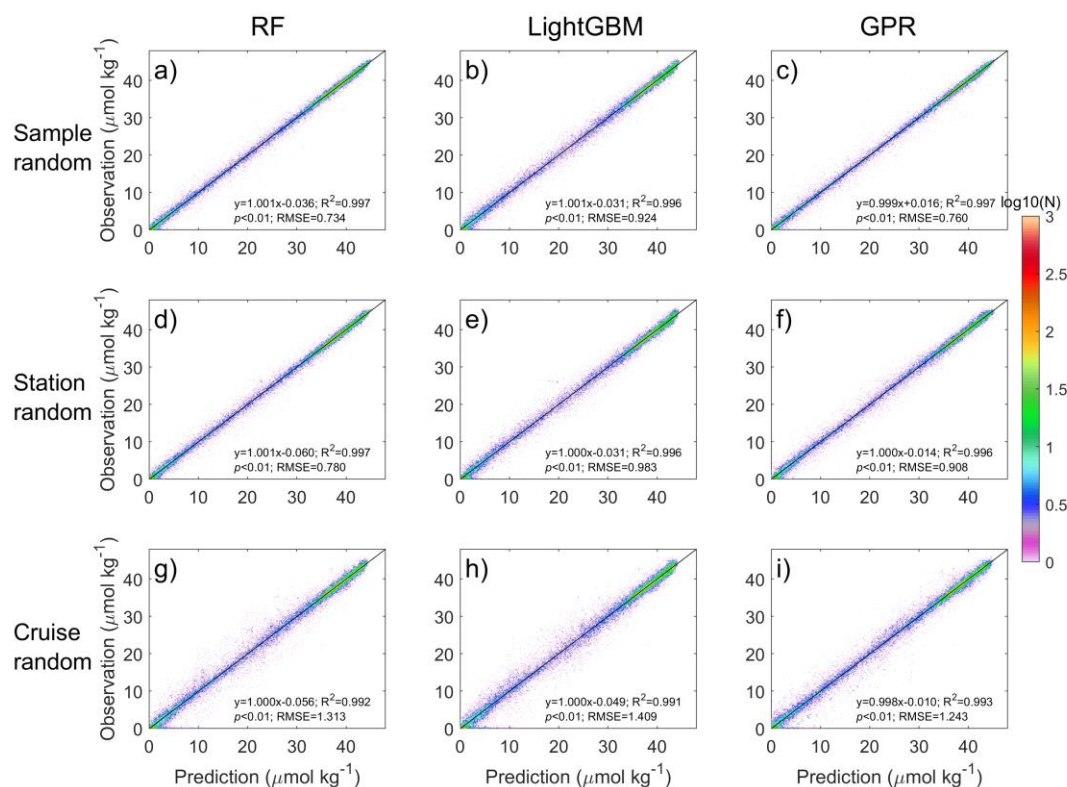
342

343 **Figure 6.** Schematic of the error estimation procedure. a) Error estimation based on
 344 three types of data selection strategy; b) assessing temporal error evolution by
 345 excluding the data at Station ALOHA; c) examining the models' reconstruction error
 346 using the hydrographic and nutrient data before 1970. The T and S denote the potential
 347 temperature and salinity, respectively.

348

349 The validation results for reconstructed NO_x^- versus observations under the first three
 350 data-selection strategies are shown in Fig. 7. RF and GPR exhibited nearly identical
 351 performance, with regression slopes of 0.992–0.998, $R^2 > 0.992$, and Root Mean
 352 Squared Errors (RMSEs) between 0.734 and 1.313 $\mu\text{mol kg}^{-1}$ (Fig. 7a, c, d, f, g, i).
 353 LightGBM showed slightly lower accuracy (slope: 0.991–0.995; R^2 : 0.991–0.996;
 354 RMSEs: 0.780–1.419 $\mu\text{mol kg}^{-1}$) (Fig. 7b, e, h). Across different data-selection
 355 strategies, sample-random (Error 1) yielded the lowest errors (RMSEs: 0.734–0.983
 356 $\mu\text{mol kg}^{-1}$) (Fig. 7a–c), station-random (Error 2) was intermediate (RMSEs: 0.908–
 357 1.313 $\mu\text{mol kg}^{-1}$) (Fig. 7d–f), and cruise-random (Error 3) produced the highest errors

358 (RMSEs: 1.243–1.424 $\mu\text{mol kg}^{-1}$) (Fig. 7; Table 3). This gradient in error estimates
 359 underscores the necessity of employing different data-selection strategies for a
 360 comprehensive error assessment. The high slopes and R^2 values (>0.99) achieved across
 361 all algorithms and data-selection strategies confirmed the robustness of the nutrient
 362 reconstructions.



363
 364 **Figure 7.** Validating the reconstructed NO_x^- concentrations using leave-one-out cross-
 365 validation with different data selection strategies and machine learning methods. Plots
 366 shown in row 1 correspond to the sample random strategy (a-c), row 2 correspond to
 367 the station random strategy (d-e), and row 3 correspond to the cruise random
 368 strategy (g-i). Plots shown in column 1 correspond to the Random Forest (RF; a, d, and
 369 g), column 2 correspond to the LightGBM (b, e, and h), and column 3 correspond to
 370 the Gaussian Process Regression (GPR; c, f, and i). The black lines and text show the
 371 fitted linear regressions, regression equations, coefficient of determination (R^2), p -
 372 values, and Root Mean Squared Errors (RMSEs). The color represents the data density
 373 (N , number of observations). Note that a logarithmic scale is applied to N .

374

375 Reconstruction errors for NO_2^- , DIP, and $\text{Si}(\text{OH})_4$ are summarized in Figs. S1–S3
 376 and Table 3. Across methods, the RMSEs were below 0.079 $\mu\text{mol kg}^{-1}$ for NO_2^- , 0.089

377 $\mu\text{mol kg}^{-1}$ for DIP, and $3.07 \mu\text{mol kg}^{-1}$ for $\text{Si}(\text{OH})_4$. DIP and $\text{Si}(\text{OH})_4$ exhibited similar
 378 error trends: RMSEs increased from sample-random to station-random to cruise-
 379 random selection. In contrast, NO_2^- reconstruction exhibited lower accuracy than NO_x^- ,
 380 DIP, and $\text{Si}(\text{OH})_4$, with regression slopes of 0.48–0.68 and R^2 values of 0.32–0.72. RF
 381 and LightGBM outperformed GPR for NO_2^- . The poorer NO_2^- performance likely
 382 reflects its generally low concentrations (mostly $<0.5 \mu\text{mol kg}^{-1}$) and high biological
 383 variability. Thus, we highlight NO_2^- as a high-uncertainty reconstruction.

384 Table 3 The Root Mean Squared Errors of nutrient reconstruction from different error
 385 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

386

387 Understanding the spatiotemporal structure of reconstruction errors is also important
 388 for assessing the models' reconstruction applicability. As shown in Figs. S4-S7, the
 389 reconstruction errors of NO_3^- , DIP, and $\text{Si}(\text{OH})_4$ are generally small in the surface layer,
 390 increase with depth to maxima at the nutricline, and then decrease to low values in deep
 391 layers. However, the random errors associated with individual cruise observations for
 392 $\text{Si}(\text{OH})_4$ display no evident vertical pattern. Horizontally, we paid particular attention
 393 to surface waters due to their greatest concentration gradients. The horizontal
 394 distribution shows that the errors are small in the western NPSG (a nutrient-depleted
 395 region) but are large in the subarctic gyre and close to the equatorial regions (nutrient-
 396 replete regions; Figs. S8-S11). Here, we particularly examined the nutrient
 397 reconstruction errors in the oligotrophic NPSG. The oligotrophic regimes are defined
 398 as regions where NO_3^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$ concentrations are <0.2 , <0.2 , <0.2 ,
 399 and $<5.0 \mu\text{mol kg}^{-1}$, respectively. As shown in Table 4, the reconstruction errors in these
 400 regimes are <0.574 , <0.056 , <0.084 , and $<1.88 \mu\text{mol kg}^{-1}$ for NO_3^- , NO_2^- , DIP, and

401 Si(OH)₄, respectively, which are evidently lower than the overall RMSEs for the entire
 402 North Pacific (Table 3). Among these models, the RF generally performs the best
 403 compared to the others. This confirms that absolute errors decrease in oligotrophic
 404 regimes. Since the number of summer observations is up to three times greater than that
 405 in winter and spring, we further examined the seasonal variation of errors. Overall, no
 406 evident seasonal variations are displayed. Only in the case of random cruise selection
 407 was the NO₃⁻ error shown to be greater in spring (March to May) than in other seasons
 408 (Fig. S12). For other cases and nutrients, seasonal variation in error was not evident.
 409 On a decadal timescale, the reconstruction errors display a slight decreasing trend,
 410 particularly for DIP, from 1973 to 2020 (Fig. S13), implying that the errors might be
 411 smaller in recent decades than in previous ones.

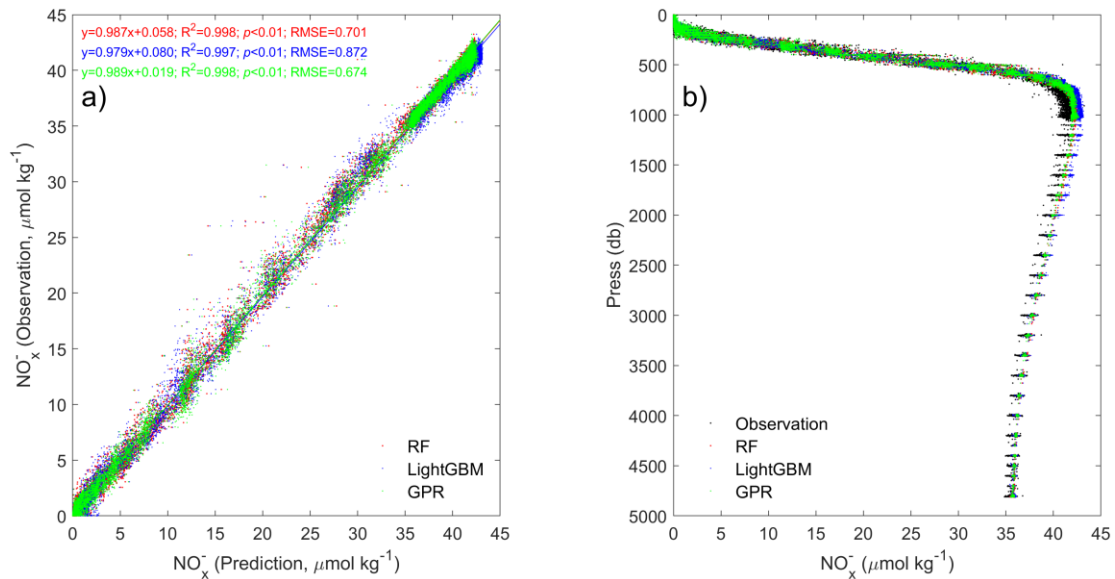
412 Table 4 The Root Mean Squared Errors of nutrient reconstruction from different error
 413 evaluation strategies in surface oligotrophic regimes (unit: $\mu\text{mol kg}^{-1}$).

Data selection strategy	NO _x ⁻			NO ₂ ⁻			DIP			Si(OH) ₄		
	RF	Light GBM	GPR	RF	Light GBM	GPR	RF	Light GBM	GPR	RF	Light GBM	GP R
Sample random	0.290	0.567	0.444	0.018	0.035	0.048	0.028	0.042	0.039	1.19	0.90	1.30
Station random	0.303	0.457	0.474	0.030	0.030	0.043	0.036	0.045	0.043	1.24	1.51	1.51
Cruise random	0.378	0.457	0.574	0.030	0.029	0.056	0.075	0.077	0.084	1.85	1.88	1.75

414

415 A fourth validation step assessed the model's temporal performance at Station
 416 ALOHA (Error 4; Fig. 6b). To test this, we withheld all observations from ALOHA
 417 (which, since 1988, represent 8.52%, 8.45%, and 8.11% of the total Si(OH)₄, NO_x⁻, and
 418 DIP records, respectively) from model training. We then reconstructed nutrient
 419 concentrations using space, time, and water-type predictors at Station ALOHA. NO₂⁻
 420 was excluded due to insufficient observations. For NO_x⁻, the regression slopes between
 421 reconstruction and observations were 0.99, 0.98, and 0.99, with RMSEs of 0.701, 0.842,
 422 and 0.674 $\mu\text{mol kg}^{-1}$ for RF, LightGBM, and GPR, respectively; R² values exceeded
 423 0.997 for all models (Fig. 8a). RF and GPR slightly outperformed LightGBM. All
 424 models accurately reproduced the NO_x⁻ profiles (Fig. 8b). The reconstruction errors for
 425 DIP were 0.066, 0.079, and 0.064 $\mu\text{mol kg}^{-1}$ for RF, LightGBM, and GPR, respectively.

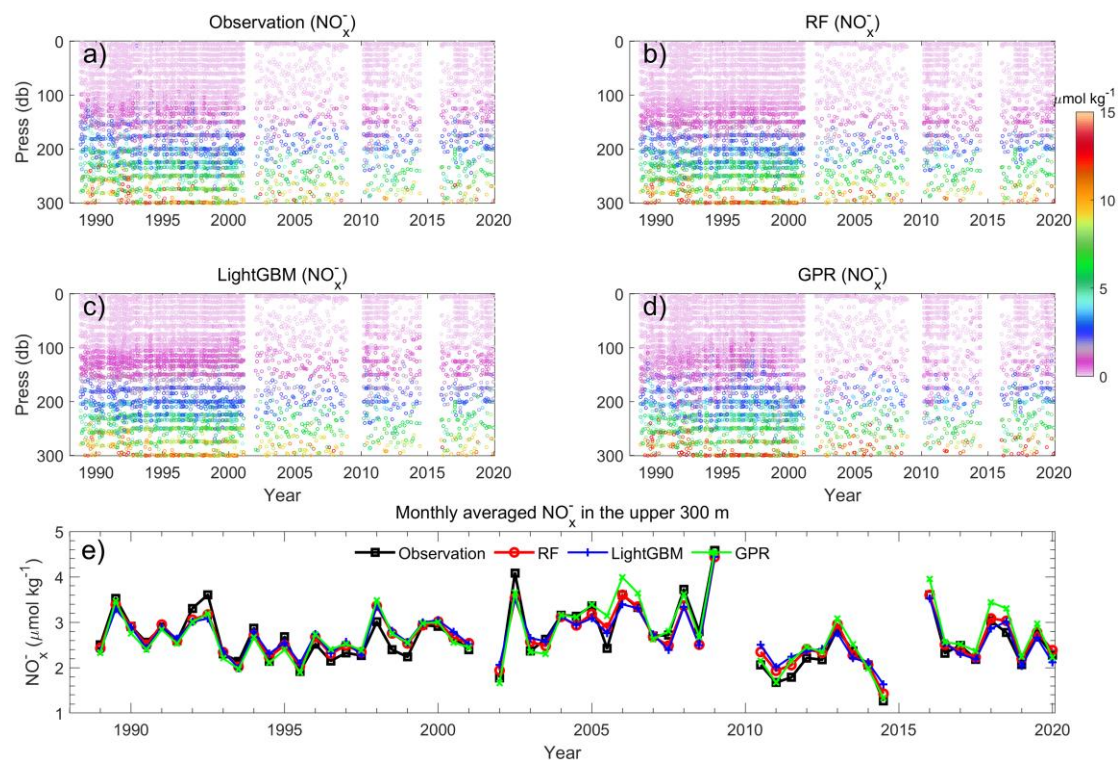
426 The corresponding errors for $\text{Si}(\text{OH})_4$ were 2.13, 2.48, and 2.32 $\mu\text{mol kg}^{-1}$ (Table 3,
427 Figs. S14–S15).



428
429 **Figure 8.** Validating the reconstructed nutrient concentrations at Station ALOHA. a)
430 Reconstructed $\text{NO}_3^- + \text{NO}_2^-$ (NO_x^-) vs. observations: Random Forest (RF; red dots),
431 LightGBM (blue dots), and Gaussian Process Regression (GPR; green dots). b) Profiles
432 of observed (black dots) and reconstructed NO_x^- from RF (red dots), LightGBM (blue
433 dots), and GPR (green dots).

434

435 Since the variations of nutrients primarily occur in the upper water column, we
436 focused on the nutrient reconstruction in the upper 300 m at Station ALOHA. Overall,
437 the models reproduced the profiles of NO_x^- from 1988 to 2021 well (Fig. 9a-d). The
438 reconstruction errors were low at the surface and increased with depth, with most of the
439 values $<3.0 \mu\text{mol kg}^{-1}$ (Fig. S16a-d). To evaluate models' ability to reconstruct nutrient
440 variations in time, the nutrient concentrations were averaged monthly over the upper
441 300 m. As compared to observations, RF, LightGBM, and GPR all well reconstructed
442 the interannual variations of NO_x^- with most of the absolute errors $<0.5 \mu\text{mol kg}^{-1}$ (Figs.
443 9e and S16e) at Station ALOHA. Similarly, the validation of DIP and $\text{Si}(\text{OH})_4$ are
444 shown in Figs. S17-S20.



445

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

451

452 A fifth validation step evaluates the models' reconstruction for the period before 1970
 453 (Error 5; Fig. 6c). This is necessary because the training data (CCHDO) spans 1973–
 454 2022, while the reconstructions are extrapolated back to 1895. We argue that this
 455 extrapolation should be reasonable because the variations of temperature-salinity-
 456 nutrient relationships in the ocean's interior might be small over the past century,
 457 providing a basis for temporal extrapolation. First, the residence time of nitrogen in
 458 deep and intermediate waters can be up to 2000 years in the North Pacific. Consequently,
 459 the imprint of centennial-scale change on nutrient inventories is attenuated. Second, the
 460 long-term variations of nutrient concentrations are not evident within our core training
 461 period (1973–2022; Figs. 9e and 17). Finally, the mean nutrient profiles derived from
 462 the 1920-1970 and 1973-2022 periods are not evidently different in the central North
 463 Pacific (Fig. S21). Therefore, while the North Pacific may experience long-term

464 variability, it might be masked by the reconstruction error, and the use of hydrographic
465 properties as predictors for nutrients is justified for historical reconstructions.

466 However, when assessing the reconstruction errors before 1970, we first consider
467 data quality issues. Prior to the standardization of modern oceanographic methods,
468 nutrient measurements—particularly from earlier decades—were subject to greater
469 analytical errors, inconsistent sampling protocols, and varied determination techniques.
470 The data quality concern is evident in the sporadic and sometimes physically
471 implausible deep nutrient profiles found in WOD for that era (Fig. S22). This is also
472 the primary reason that nutrient data pre-1973 collected from sources like the OSD from
473 WOD were not incorporated into model training. To evaluate data quality in earlier
474 decades, we selected five specific years with more abundant observations: 1929, 1947,
475 1953, 1958, and 1966 (Fig. S23). After applying the same quality-control criteria
476 outlined in Section 3.1, we used the historical hydrographic data (temperature and
477 salinity) from those years to predict nutrient concentrations. A total of 52,277, 119,137,
478 284,472, and 193,339 data points were collected for NO_3^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$,
479 respectively, after QC. The comparison between these predictions and the quality-
480 controlled observations yields the prediction errors for the pre-1970 period (Fig. 6c).
481 The RMSEs from different models suggested values <5.7 , <0.40 , and $<22.9 \mu\text{mol kg}^{-1}$
482 for NO_3^- , DIP, and $\text{Si}(\text{OH})_4$, respectively (Figs. S24–S26), which are much larger than
483 the corresponding errors for the period after 1970. We recommend that these values be
484 considered a conservative estimate of the upper error bound, as they incorporate both
485 nutrient observations and prediction errors. In addition, the hydrographic data are also
486 less reliable in the earlier period. Thus, we acknowledge that reconstruction errors are
487 likely higher for the pre-1973 period, and the error estimated here should be considered
488 a "best estimate" with quantified uncertainties, and encourage users to consider these
489 error bounds when applying the dataset to early twentieth-century conditions.

490

491 **3.2 Reconstructed nutrients**

492 The final reconstructed nutrient dataset aligns with the spatiotemporal coverage of
493 the quality-controlled WOD hydrographic dataset, comprising 472,652,680 data points
494 for each nutrient (NO_x^- , NO_2^- , DIP, and $\text{Si}(\text{OH})_4$) from 1,920,634 stations across 35,744
495 cruises, spanning from 1895 to 2024 (Table 2). Most data points are located above 2,000
496 m, with fewer observations at greater depths due to hydrographic platform limitations.

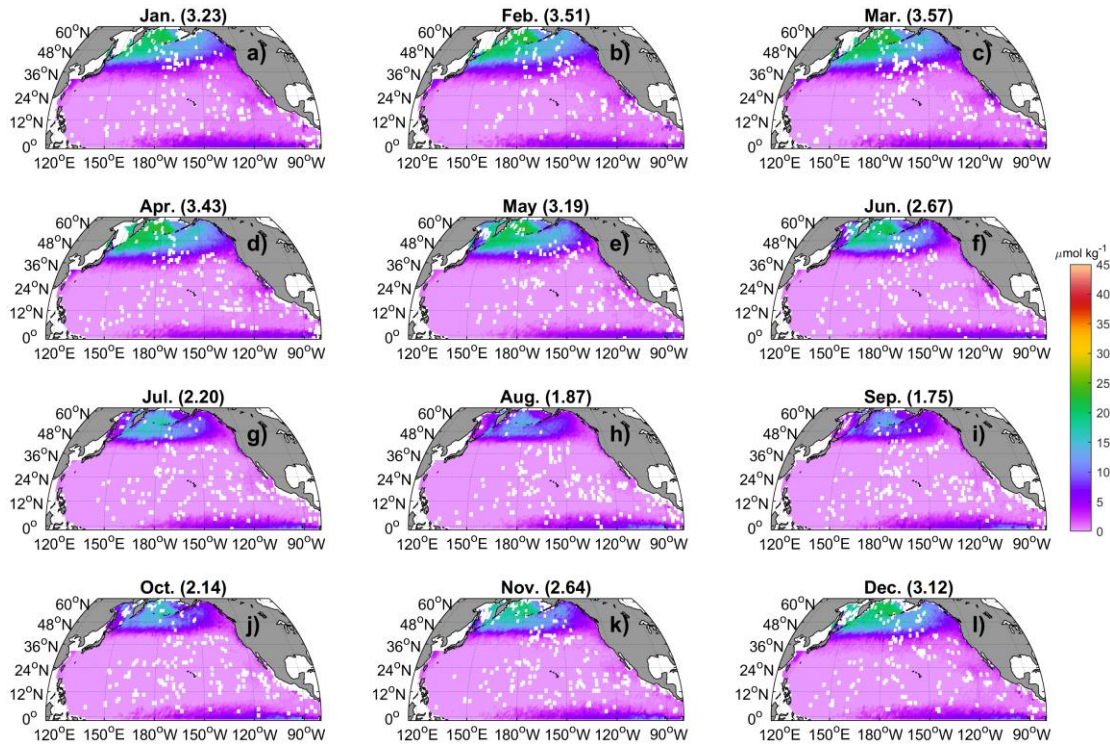
497 It is important to clarify the nature of the reconstructed dataset, which is
498 fundamentally different from gridded products. This product provides nutrient
499 concentrations linked to each hydrographic observations: nutrient values are
500 reconstructed precisely at the locations, depths, and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

505 **3.3 Climatology of nutrient distributions**

506 To evaluate the reliability of our product, we binned and averaged the predicted
507 nutrients within $1^\circ \times 1^\circ$ grid cells for each month to produce a monthly climatology. This
508 climatology represents a mean field that depends heavily on the spatiotemporal
509 distribution of the underlying data and may be influenced by uneven data sampling.
510 This reconstructed climatology was compared with the World Ocean Atlas 2023
511 (WOA23), which is derived from quality-controlled and objectively analyzed
512 observational data. Since the large-scale patterns of NO_3^- , DIP, and $\text{Si}(\text{OH})_4$ are similar
513 among different models (Figs. 10–13, S27–S36), we focus on NO_3^- reconstructed by
514 the RF model in this section unless stated otherwise.

515 Figs. 10–13 present the monthly climatology of NO_x^- at 5 m, 100 m, 500 m, and
516 1,000 m in the North Pacific. At 5 m, the reconstructed NO_x^- accurately captures the
517 established spatial patterns, with elevated concentrations in the subpolar gyre, Bering
518 Sea, and equatorial regions, and depleted concentrations in the NPSG (Fig. 10).
519 Seasonally, the basin-averaged surface NO_x^- concentrations display the highest value
520 of $3.50 \mu\text{mol kg}^{-1}$ in March, in contrast to the lowest value of $1.82 \mu\text{mol kg}^{-1}$ in

521 September. These results agree with Yasunaka et al. (2014, 2021), who, using extensive
 522 surface nutrient observations (up to 14,000 for nitrate) in the North Pacific, reported
 523 similar spatial and seasonal patterns.

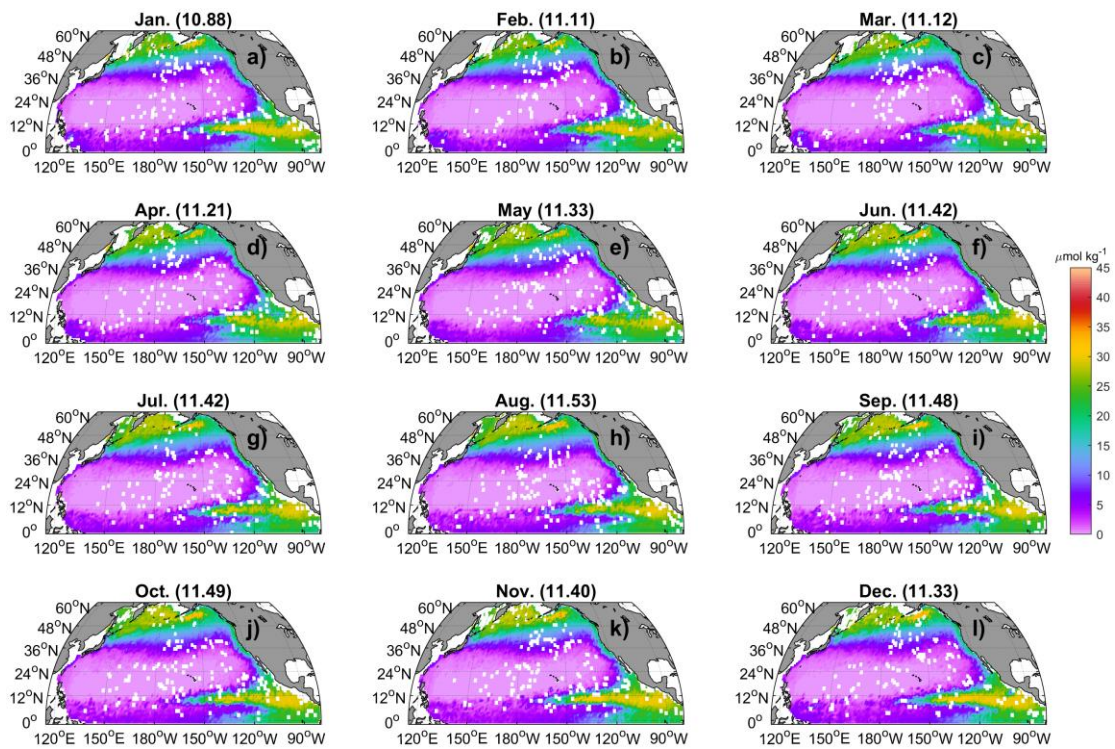


524 **Figure 10.** The monthly climatology of NO_x^- at 5 m in the North Pacific. Data are
 525 binned and averaged within $1 \times 1^\circ$ grid cells. The values in the title represent the spatial
 526 mean values.
 527

528

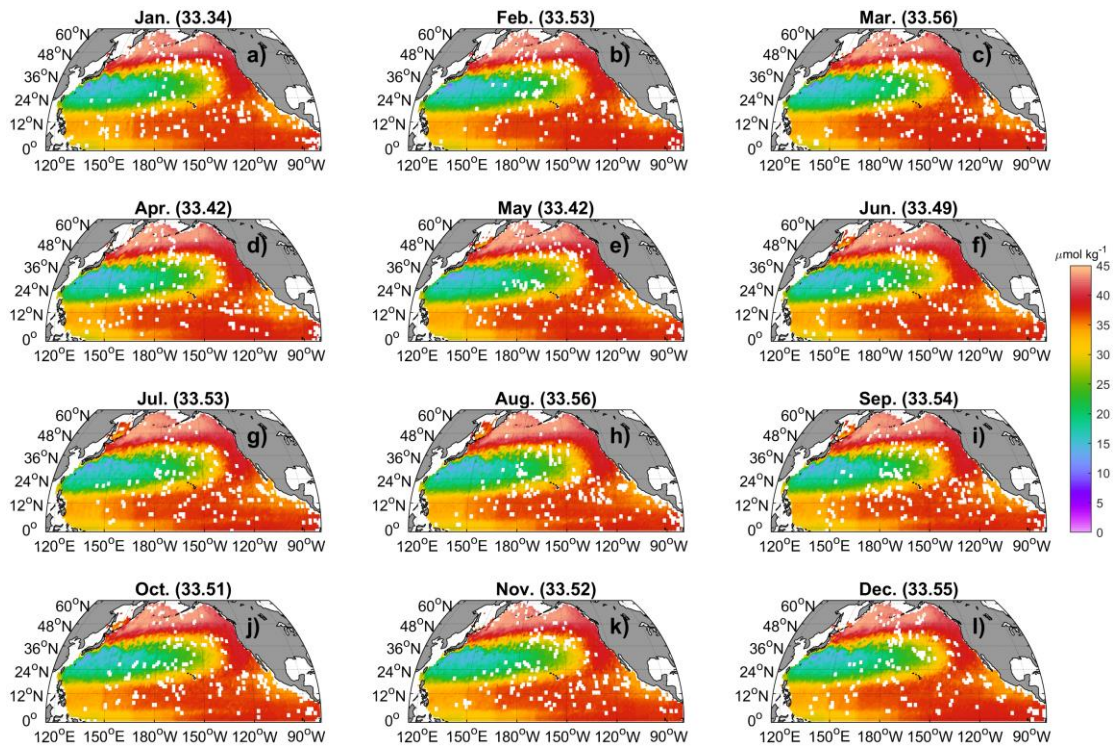
529 At 100 m, NO_x^- concentrations are elevated particularly in the subarctic gyre, north
 530 of the Equator, and the eastern North Pacific, while the central regions, particularly the
 531 NPSG, exhibit lower values. At 500 m, NO_x^- concentrations display patterns similar to
 532 those at 100 m, except that the NO_x^- concentrations in the western NPSG are evidently
 533 lower than those in other regions (Fig. 12). At 1000 m, concentrations in the
 534 southwestern North Pacific Ocean are markedly lower than those in other regions (Fig.
 535 13). Below 100 m depth, seasonal variability in NO_x^- is minimal (Figs. 11–13). These
 536 results display patterns similar to WOA23 (Figs. S36–S44). The differences between
 537 the averaged values of these two climatologies are generally $<0.7 \mu\text{mol kg}^{-1}$ at the
 538 surface and $<1.5 \mu\text{mol kg}^{-1}$ at 100 m and 500 m. The maximum differences are found

539 in July at a depth of 500 m (Figs. 13g and S38g). In that month and layer, WOA23
 540 shows a notably low mean NO_3^- value ($31.94 \mu\text{mol kg}^{-1}$) compared to its values in other
 541 months (33.15 to $34.64 \mu\text{mol kg}^{-1}$; Fig. S38) and compared to our climatology (33.34
 542 to $33.56 \mu\text{mol kg}^{-1}$; Fig. 13). This discrepancy arises because the WOA23 climatology
 543 for July features a pronounced low- NO_3^- patch (down to $20 \mu\text{mol kg}^{-1}$) within the
 544 eastern subarctic gyre, surrounded by waters with concentrations of $>35 \mu\text{mol kg}^{-1}$ (Fig.
 545 S38g). These regional differences are clearly visible in the difference maps between the
 546 two products (Figs. S45–S47). Generally, our reconstructions capture finer spatial detail,
 547 exhibit less oversmoothing, and avoid artificial “bull’s-eye” patterns.



548
 549 **Figure 11.** The monthly climatology of NO_x^- at 100 m in the North Pacific. Data are
 550 binned and averaged within $1 \times 1^\circ$ grid cells. The values in the title represent the spatial
 551 mean values.

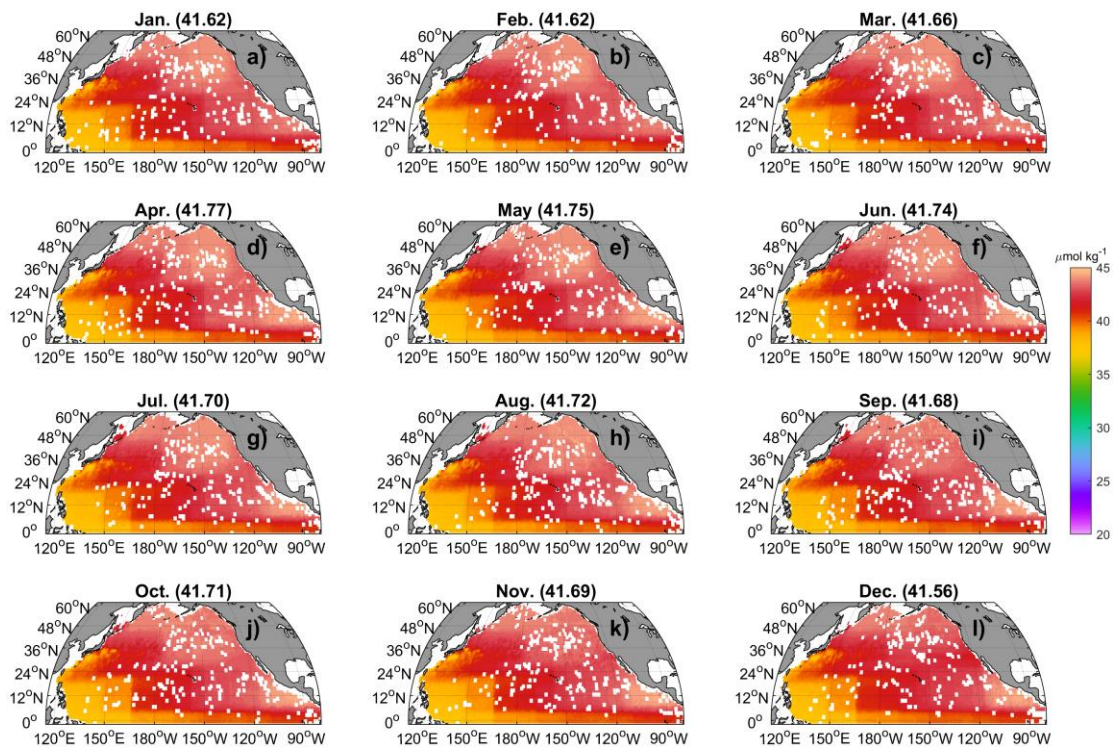
552



553

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

557

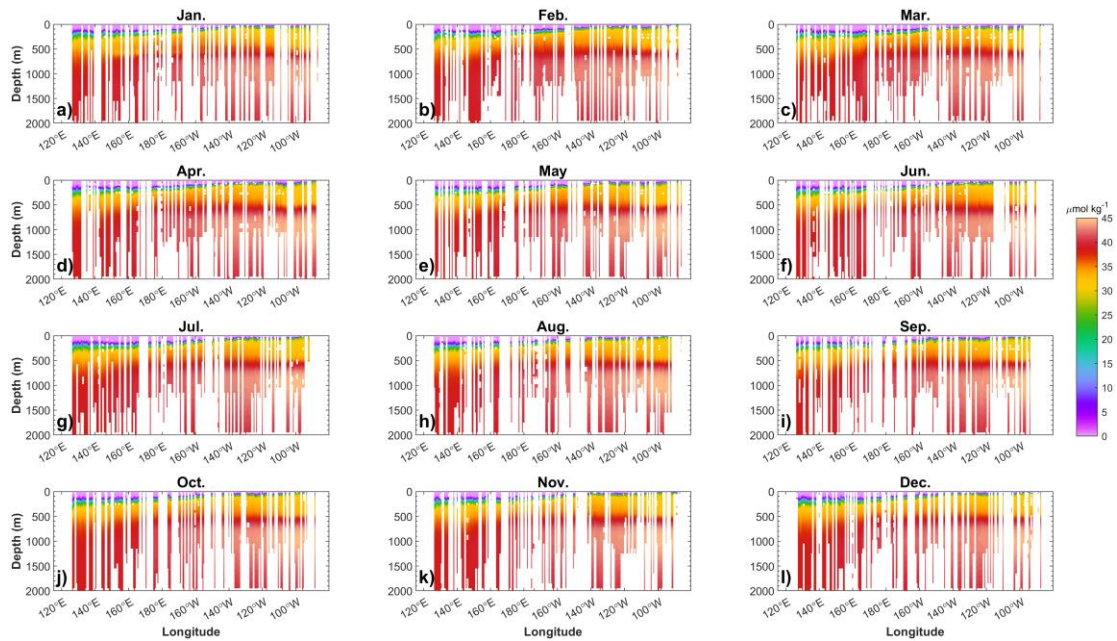


558

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

562

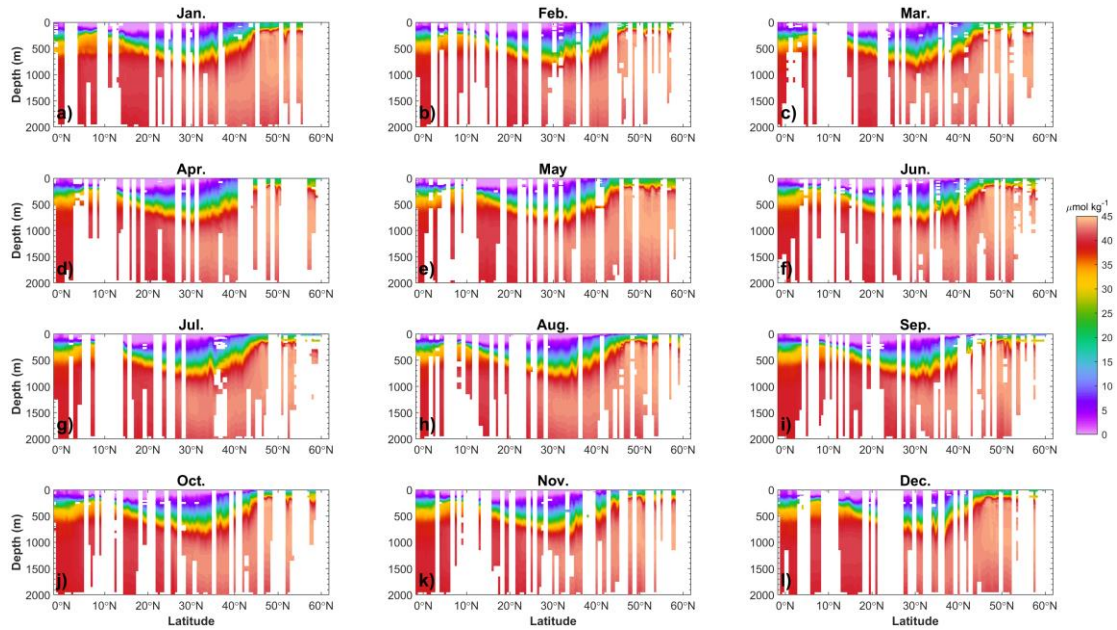
563 Sectional distributions of NO_x^- in the upper 2000 m along 10° N and 180° E were
 564 used as examples to illustrate the vertical profile distributions of nutrients within the
 565 North Pacific. At 10° N, NO_x^- concentrations increase from $\sim 0.0 \mu\text{mol kg}^{-1}$ at the
 566 surface to $\sim 45.0 \mu\text{mol kg}^{-1}$ at ~ 1000 m, followed by a decrease to $\sim 38.0 \mu\text{mol kg}^{-1}$ at
 567 2000 m. NO_x^- concentrations increase from west to the east in the North Pacific in the
 568 upper 300 m (Fig. 14). At 180° E, in the upper 500 m, meridional NO_x^- concentrations
 569 increase from the equator to the North Equatorial Current ($\sim 10^\circ$ N), decline within the
 570 subtropical gyre, and then increase toward the subarctic region (Fig. 15). Generally,
 571 seasonal differences of NO_x^- concentrations along both sections are not evident.



572

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

575



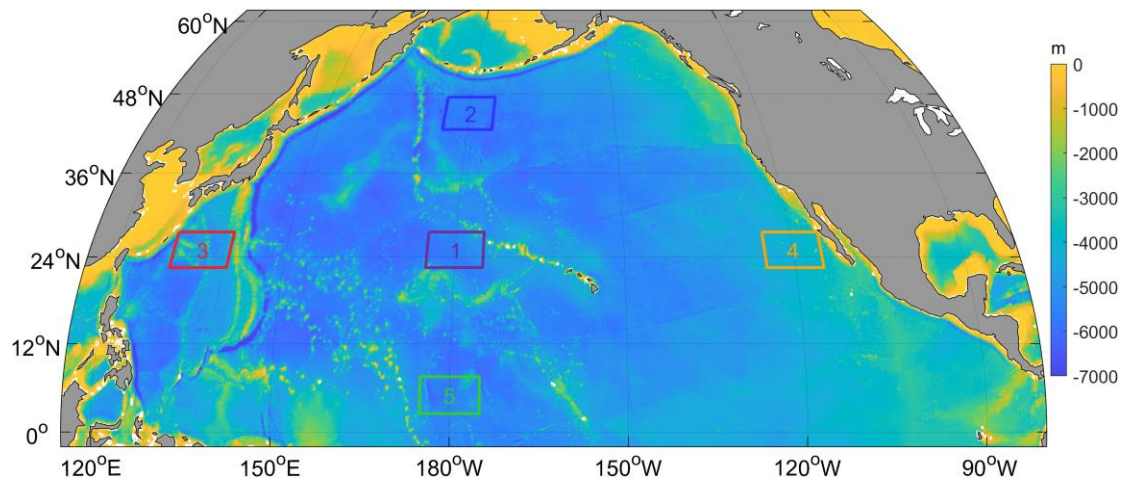
576

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

579

580 3.4 Long-term variations of nutrients

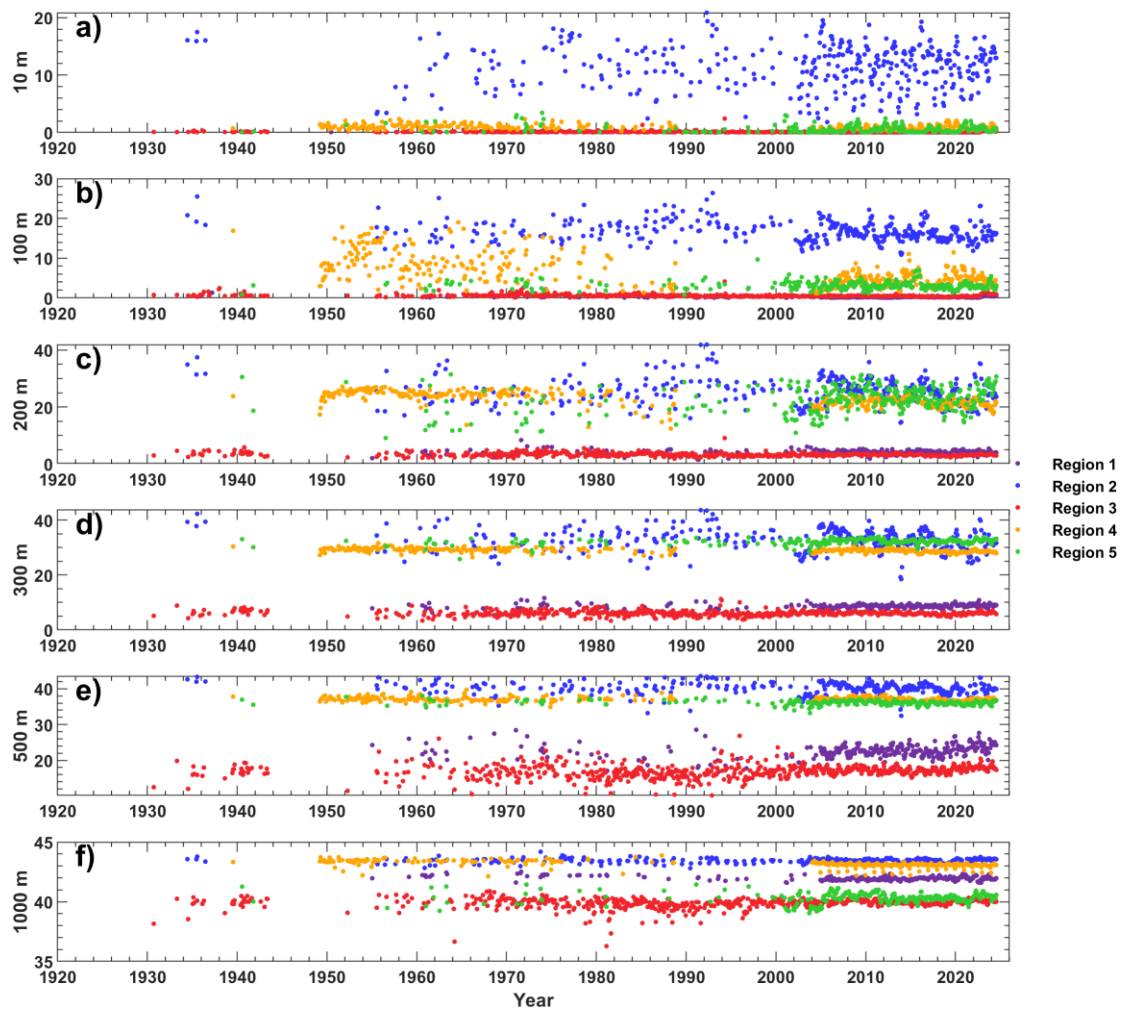
581 We present an initial analysis of long-term nutrient changes by examining five
 582 representative regions in the North Pacific, covering the subarctic gyre, the subtropical
 583 gyre, and equatorial areas (Fig. 16). The data are binned by region, month, and depth
 584 (10 m, 100 m, 200 m, 300 m, 500 m, and 1000 m) for regions 1–5. As shown in Fig. 17,
 585 these time series reveal notable interannual fluctuations of NO_3^- (with 2–5-year
 586 oscillations), providing a first-order view of low-frequency variability captured by the
 587 reconstruction. However, no evident long-term trend is found for nutrients. DIP and
 588 $\text{Si}(\text{OH})_4$ display patterns similar to NO_3^- (Figs. S48–S49). In contrast, at depths of 200
 589 m and 300 m, NO_2^- displays an increasing trend in the central NPSG and a decreasing
 590 trend in the eastern NPSG during the 1970–2005 period (Fig. S50). More sophisticated
 591 trend analyses and basin-scale integrations are promising avenues for future work based
 592 on this newly reconstructed dataset.



593

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

596



597

598 **Figure 17.** Time series of reconstructed NO_3^- concentrations at 10 m (a), 100 m (b),
599 200 m (c), 300 m (d), 500 m (e), and 1000 m (f) for regions 1–5 (see Fig. 16). Data
600 were binned by depth and region and then averaged by month.

601

602 **4 Data availability**

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

610

611 **5 Conclusion**

612 In this study, we applied rigorous quality control procedures to clean hydrographic
613 and nutrient observations from CCHDO and WOD datasets. The cleaned CCHDO data
614 were then used to train three machine-learning models to relate nutrient concentrations
615 to spatial, temporal, and water-mass predictors. The models were applied to reconstruct
616 nutrient concentrations from hydrographic observations collected from WOD, most
617 of which lack direct nutrient measurements. We assessed the model performance using
618 four data-partition strategies, and found that all models reproduced held-out data with
619 low RMSEs. RF and GPR slightly outperformed LightGBM. The application of these
620 models to WOD hydrography yielded 472,652,680 reconstructed nutrient
621 concentrations across 1,920,634 stations and 35,744 cruises, spanning from 1895 to
622 2024. This represents a 2,127– to 2,393-fold increase compared to the original volume
623 of CCHDO nutrient data. The reconstruction captured the spatial, seasonal, and
624 interannual variations of water column nutrients in the North Pacific Ocean well.
625 Compared to the WOA23 climatology, the reconstruction-based nutrient climatology
626 exhibited more realistic spatial structures than WOA23. This high-quality and high-

627 resolution nutrient dataset adds historical nutrient estimation for locations and times
628 with solely hydrographic measurements. Additional potential applications of this
629 dataset include: 1) investigating nutrient transport and budget in the North Pacific; 2)
630 spinning up and validating ocean biogeochemical models; 3) assessing long-term
631 nutrient trends driven by anthropogenic forcing and climate change; 4) investigating
632 nutrient stoichiometric changes and their ecological impacts under climate variability.
633 Collectively, this resource facilitates advanced studies on marine biogeochemical
634 cycles, ecosystem dynamics, and climate-nutrient interactions.

635

636 **Author contributions**

637 CD and XL designed the study and dataset. CD, SK, MD, ZC, DS, and XL conceived
638 the project and secured the funding. CD, NZ, QL, HW and XL collected and processed
639 the data, developed the code, and performed the analysis. SK, MD, ZC, and DS
640 provided methodological guidance and advice. CD and NZ wrote the original draft. All
641 authors reviewed and edited the manuscript.

642

643 **Competing interests**

644 The corresponding author has declared that none of the authors has any competing
645 interests.

646

647 **Acknowledgements**

648 This study was funded by the National Key R&D Program of China
649 (Grant 2023YFF0805001), This study was funded by the National Natural Science
650 Foundation of China (Grants 42494885, 42576215, 42494881, 42276034),
651 Fundamental and Interdisciplinary Disciplines Breakthrough Plan of the Ministry of
652 Education of China (Grant JYB2025XDXM801), Innovational Fund for Scientific and
653 Technological Personnel of Hainan Province (Grant KJRC2023B04), and Natural
654 Science Foundation of Hainan Province (Grant 624MS037). We thank the CCHDO
655 (<https://cchdo.ucsd.edu/>) and the WOD (<https://www.ncei.noaa.gov/products/world->

656 ocean-database) for providing the data used in this study. Special thanks are owed to all
657 scientists involved in data collection, analysis, and management for these programs.

658

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

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

663

664 **References**

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

668 Ascani, F., Richards, K. J., Firing, E., Grant, S., Johnson, K. S., Jia, Y., Lukas, R., and
669 Karl, D. M.: Physical and biological controls of nitrate concentrations in the upper
670 subtropical North Pacific Ocean, *Deep-Sea Res. Pt. II*, 93, 119–134, 2013.

671 Barone, B., Church, M. J., Dugenne, M., Hawco, N. J., Jahn, O., White, A. E., John, S.
672 G., Follows, M. J., DeLong, E. F., and Karl, D. M.: Biogeochemical dynamics in
673 adjacent mesoscale eddies of opposite polarity, *Global Biogeochem. Cy.*, 36,
674 e2021GB007115, 2022.

675 Benitez-Nelson, C. R., Bidigare, R. R., Dickey, T. D., Landry, M. R., Leonard, C. L.,
676 Brown, S. L., Nencioli, F., Rii, Y. M., Maiti, K., Becker, J. W., Bibby, T. S., Black,
677 W., Cai, W. J., Carlson, C. A., Chen, F., Kuwahara, V. S., Mahaffey, C., McAndrew,
678 P. M., Quay, P. D., Rappé, M. S., Selph, K. E., Simmons, M. P., and Yang, E. J.:
679 Mesoscale Eddies Drive Increased Silica Export in the Subtropical Pacific Ocean,
680 *Science*, 316, 1017–1021, 2007.

681 Bidigare, R. R., Chai, F., Landry, M. R., Lukas, R., Hannides, C. C. S., Christensen, S.
682 J., Karl, D. M., Shi, L., and Chao, Y.: Subtropical ocean ecosystem structure

683 changes forced by North Pacific climate variations, *J. Plankton Res.*, 31, 1131–
684 1139, 2009.

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

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

690 Browning, T. J., Liu, X., Zhang, R., Wen, Z., Liu, J., Zhou, Y., Xu, F., Cai, Y., Zhou, K.,
691 Cao, Z., Zhu, Y., Shi, D., Achterberg, E. P., and Dai, M.: Nutrient co-limitation in
692 the subtropical Northwest Pacific, *Limnol. Oceanogr. Lett.*, 7, 52–61, 2021.

693 Chelton, D. B., Schlax, M. G., Samelson, R. M., and de Szoeko, R. A.: Global
694 observations of large oceanic eddies, *Geophys. Res. Lett.*, 34, L15606, 2007.

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

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

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

703 Dai, M., Luo, Y., Achterberg, E. P., Browning, T. J., Cai, Y., Cao, Z., Chai, F., Chen, B.,
704 Church, M. J., Ci, D., Du, C., Gao, K., Guo, X., Hu, Z., Kao, S., Laws, E. A., Lee,
705 Z., Lin, H., Liu, Q., Liu, X., Luo, W., Meng, F., Shang, S., Shi, D., Saito, H., Song,
706 L., Wan, X. S., Wang, Y., Wang, W.-L., Wen, Z., Xiu, P., Zhang, J., Zhang, R., and
707 Zhou, K.: Upper Ocean biogeochemistry of the oligotrophic North Pacific
708 subtropical gyre: From nutrient sources to carbon export, *Rev. Geophys.*, 61,
709 e2022RG000800, 2023.

710 Du, C., Zheng, N., Kao, S.-J., Dai, M., Cao, Z., Shi, D., Li, Q., Wang, H., and Li, X.:
711 Validated temperature and salinity data, and reconstructed nutrient concentrations

712 in the North Pacific (1895 – 2024). Zenodo, <https://zenodo.org/records/17451417>,
713 2025.

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

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

718 Dong, L., Qi, J., Yin, B., Zhi, H., Li, D., Yang, S., Wang, W., Cai, H., and Xie, B.:
719 Reconstruction of subsurface salinity structure in the South China Sea using
720 satellite observations: a LightGBM-Based Deep forest method, *Remote Sens.*, 14,
721 3494, 2022.

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

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

728 Eugster, O. and Gruber, N.: A probabilistic estimate of global marine N-fixation and
729 denitrification, *Glob. Biogeochem. Cycles*, 26, GB4013, 2012.

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

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

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

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

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

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

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

752 Kamykowski, D., Zentara, S.-J., Morrison, J. M., and Switzer, A. C.: Dynamic global
753 patterns of nitrate, phosphate, silicate, and iron availability and phytoplankton
754 community composition from remote sensing data, *Global Biogeochem. Cy.*, 16,
755 1077, 2002.

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

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

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

762 Karl, D. M., Letelier, R. M., Bidigare, R. R., Björkman, K. M., Church, M. J., Dore, J.
763 E., and White, A. E.: Seasonal-to-decadal scale variability in primary production
764 and particulate matter export at Station ALOHA, *Prog. Oceanogr.*, 195, 102563,
765 2021.

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

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

771 Liaw, A. and Wiener, M.: Classification and regression by randomForest, *R News*, 2,
772 18–22, 2002.

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

775 Liu, H., Lin, L., Wang, Y., Du, L., Wang, S., Zhou, P., Yu, Y., Gong, X., and Lu, X.:
776 Reconstruction of Monthly Surface Nutrient Concentrations in the Yellow and
777 Bohai Seas from 2003–2019 Using Machine Learning, *Remote Sens.*, 14, 5021,
778 2022.

779 Madani, N., Parazoo, N. C., Manizza, M., Chatterjee, A., Carroll, D., Menemenlis, D.,
780 Fouest, V. L., Matsuoka, A., Luis, K. M., Serra-Pompei, C., and Miller, C. E.: A
781 machine learning approach to produce a continuous Solar-Induced chlorophyll
782 fluorescence over the Arctic Ocean, *J. Geophys. Res. Machine Learn. Comput.*, 1,
783 2024.

784 Martino, M., Hamilton, D. S., Baker, A. R., Jickells, T., Bromley, T., Nojiri, Y., Quack,
785 B., and Boyd, P. W.: Western Pacific atmospheric nutrient deposition fluxes, their
786 impact on surface ocean productivity, *Glob. Biogeochem. Cycles*, 28, 712–728,
787 2014.

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

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

797 Możejko, J. and Gniot, R.: Application of Neural Networks for the Prediction of Total
798 Phosphorus Concentrations in Surface Waters, *Pol. J. Environ. Stud.*, 17, 363–368,
799 2008.

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

803 Qi, J., Yu, Y., Yao, X., Yuan, G., and Gao, H.: Dry deposition fluxes of inorganic
804 nitrogen and phosphorus in atmospheric aerosols over the Marginal Seas and
805 Northwest Pacific, *Atmos. Res.*, 245, 105076, 2020.

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

810 Sarangi, P. K., Thangaradjou, T., Kumar, A. S., and Balasubramanian, T.: Development
811 of nitrate algorithm for the southwest bay of bengal water and its implication using
812 remote sensing satellite datasets, *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*,
813 4, 983–991, 2011.

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

816 Steinhoff, T., Friedrich, T., Hartman, S. E., Oeschies, A., Wallace, D. W. R., and
817 Körtzinger, A.: Estimating mixed layer nitrate in the North Atlantic Ocean,
818 *Biogeosciences*, 7, 795–807, 2010.

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

822 Sundararaman, H. K. K. and Shanmugam, P.: Estimates of the global ocean surface
823 dissolved oxygen and macronutrients from satellite data, *Remote Sens. Environ.*,
824 311, 114243, 2024.

825 Switzer, A. C., Kamykowski, D., and Zentara, S.-J.: Mapping nitrate in the global ocean
826 using remotely sensed sea surface temperature, *J. Geophys. Res.*, 108, 345–359,
827 2003.

828 Talley, L. D., Pickard, G. L., Emery, W. J., and Swift, J. H.: *Descriptive Physical*
829 *Oceanography, An Introduction, Sixth Edition*, Academic Press, 350–362 pp.,
830 2011.

831 Wang, C., Su, B., Sun, J., Hu, X., and Liu, J.: A regional ocean database for the Coastal
832 China Sea. *Sci Data*, 12, 1550, 2025.

833 Wang, L., Xu, Z., Gong, X., Zhang, P., Hao, Z., You, J., Zhao, X., and Guo, X.:
834 Estimation of nitrate concentration and its distribution in the northwestern Pacific
835 Ocean by a deep neural network model, *Deep Sea Res. I*, 195, 104005, 2023.

836 Wang, W.-L., Moore, J. K., Martiny, A. C., and Primeau, F. W.: Convergent estimates
837 of marine nitrogen fixation, *Nature*, 566, 205–211, 2019.

838 Wang, Z., Wang, G., Guo, X., Hu, J., and Dai, M.: Reconstruction of High-Resolution
839 Sea Surface Salinity over 2003–2020 in the South China Sea Using the Machine
840 Learning Algorithm LightGBM Model, *Remote Sens.*, 14, 6147, 2022.

841 Yang, G. G., Wang, Q., Feng, J., He, L., Li, R., Lu, W., Liao, E., and Lai, Z.: Can three-
842 dimensional nitrate structure be reconstructed from surface information with
843 artificial intelligence? – A proof-of-concept study, *Sci. Total Environ.*, 924,
844 171365, 2024.

845 Yasunaka, S., Mitsudera, H., Whitney, F., and Nakaoka, S.: Nutrient and dissolved
846 inorganic carbon variability in the North Pacific, *J. Oceanogr.*, 77, 3–16, 2021.

847 Yasunaka, S., Nojiri, Y., Nakaoka, S., Ono, T., Whitney, F. A., and Telszewski, M.:
848 Mapping of sea surface nutrients in the North Pacific: Basin-wide distribution and
849 seasonal to interannual variability, *J. Geophys. Res. Oceans*, 119, 7756–7771,
850 2014.

851 Yasunaka, S., Ono, T., Nojiri, Y., Whitney, F. A., Wada, C., Murata, A., Nakaoka, S.,
852 and Hosoda, S.: Long-term variability of surface nutrient concentrations in the
853 North Pacific, *Geophys. Res. Lett.*, 43, 3389–3397, 2016.

854 Yu, X. R., Wen, Z., Jiang, R., Yang, J.-Y. T., Cao, Z., Hong, H., Zhou, Y., and Shi, D.:
855 Assessing N₂ fixation flux and its controlling factors in the (sub)tropical western
856 North Pacific through high-resolution observations, *Limnol. Oceanogr. Lett.*, 9,
857 716 – 724, 2024.

858 Zhong, A., Wang, D., Gong, F., Zhu, W., Fu, D., Zheng, Z., Huang, J., He, X., and Bai,
859 Y.: Remote sensing estimates of global sea surface nitrate: Methodology and
860 validation, *Sci. Total Environ.*, 950, 175362, 2024.

861

862