

# 1 Global biogeography of N<sub>2</sub>-fixing microbes: *nifH* amplicon database 2 and analytics workflow

3 Michael Morando<sup>1\*</sup>, Jonathan Magasin<sup>1\*</sup>, Shunyan Cheung<sup>1,2</sup>, Matthew M. Mills<sup>3</sup>, Jonathan P. Zehr<sup>1</sup>,  
4 Kendra A. Turk-Kubo<sup>1</sup>

5 <sup>1</sup>Ocean Sciences Department, University of California, Santa Cruz, Santa Cruz, 95064, United States

6 <sup>2</sup>Institute of Marine Biology and Center of Excellence for the Oceans, National Taiwan Ocean University, Keelung, Taiwan

7 <sup>3</sup>Earth System Science, Stanford University, Stanford, 94305, United States

8 \* equal contributions

9 Correspondence to: Kendra A. Turk-Kubo ([kturk@ucsc.edu](mailto:kturk@ucsc.edu))

10 **Abstract.** Marine dinitrogen (N<sub>2</sub>) fixation is a globally significant biogeochemical process carried out by a specialized group  
11 of prokaryotes (diazotrophs), yet our understanding of their ecology is constantly evolving. Although marine N<sub>2</sub> fixation is  
12 often ascribed to cyanobacterial diazotrophs, indirect evidence suggests that non-cyanobacterial diazotrophs (NCDs) might  
13 also be important. One widely used approach for understanding diazotroph diversity and biogeography is polymerase chain  
14 reaction (PCR)-amplification of a portion of the *nifH* gene, which encodes a structural component of the N<sub>2</sub>-fixing enzyme  
15 complex, nitrogenase. An array of bioinformatic tools exists to process *nifH* amplicon data, however, the lack of  
16 standardized practices has hindered cross-study comparisons. This has led to a missed opportunity to more thoroughly assess  
17 diazotroph diversity, biogeography, and their potential contributions to the marine N cycle. To address these knowledge gaps  
18 a bioinformatic workflow was designed that standardizes the processing of *nifH* amplicon datasets originating from  
19 high-throughput sequencing (HTS). Multiple datasets are efficiently and consistently processed with a specialized DADA2  
20 pipeline to identify amplicon sequence variants (ASVs). A series of customizable post-pipeline stages then detect and discard  
21 spurious *nifH* sequences and annotate the subsequent quality-filtered *nifH* ASVs using multiple reference databases and  
22 classification approaches. This newly developed workflow was used to reprocess nearly all publicly available *nifH* amplicon  
23 HTS datasets from marine studies, and to generate a comprehensive *nifH* ASV database containing 9383 ASVs aggregated  
24 from 21 studies that represent the diazotrophic populations in the global ocean. For each sample, the database includes  
25 physical and chemical metadata obtained from the Simons Collaborative Marine Atlas Project (CMAP). Here we  
26 demonstrate the utility of this database for revealing global biogeographical patterns of prominent diazotroph groups and  
27 highlight the influence of sea surface temperature. The workflow and *nifH* ASV database provide a robust framework for  
28 studying marine N<sub>2</sub> fixation and diazotrophic diversity captured by *nifH* amplicon HTS. Future datasets that target  
29 understudied ocean regions can be added easily, and users can tune parameters and studies included for their specific focus.

30 The workflow and database are available, respectively, in GitHub (<https://github.com/jdmagasin/nifH-ASV-workflow>;  
31 Morando et al., 2024c) and Figshare (<https://doi.org/10.6084/m9.figshare.23795943.v2>; Morando et al., 2024a).

## 32 1 Introduction

33 Dinitrogen ( $N_2$ ) fixation, the reduction of  $N_2$  into bioavailable  $NH_3$  is a source of new nitrogen (N) in the oceans and can  
34 support as much as 70 % of new primary production in N-limited oligotrophic gyres (Jickells et al., 2017). Over millennia,  
35  $N_2$  fixation may balance the loss of N from the marine system through denitrification and annamox (Zehr and Capone, 2020).  
36  $N_2$  fixation was thought to be performed exclusively by prokaryotes, yet it was recently demonstrated that the marine  
37 haptophyte alga, *Braarudosphaera bigelowii*, contains a cyanobacterially-derived organelle specialized for  $N_2$  fixation  
38 (Coale et al., 2024). Noting this exception, microorganisms able to fix  $N_2$  (diazotrophs), are broadly characterized into two  
39 main groups, cyanobacterial diazotrophs (those phylogenetically related to cyanobacteria) and non-cyanobacterial  
40 diazotrophs (NCDs). Historically, cyanobacterial diazotrophs have been considered the most important contributors to  
41 marine  $N_2$  fixation (Villareal, 1994; Capone et al., 2005). NCDs, first detected by Zehr et al. (1998), have since been  
42 demonstrated to be ubiquitous in pelagic marine waters, and are generally thought to be putative chemoheterotrophs with a  
43 highly diverse lineage that includes the massive phylum Proteobacteria as well as Firmicutes, Actinobacteria, and  
44 Chloroflexi (Turk-Kubo et al., 2022). However, their contribution of fixed N and their role in the global ocean is not  
45 well-understood (Moisander et al., 2017).

46  
47 Diazotrophs are often present at low abundances relative to other members of ocean microbiomes, which makes them  
48 challenging to study (Moisander et al., 2017; Benavides et al., 2021). Distinctive pigments and morphologies that enable  
49 some cyanobacterial diazotrophs to be identified by microscopy are lacking in many diazotrophs (Carpenter and Capone,  
50 1983; Carpenter and Foster, 2002), including NCDs. Furthermore, many marine diazotrophs are uncultivated, which has  
51 required the use of cultivation-independent approaches such as PCR and quantitative PCR (qPCR) (Luo et al., 2012; Shao  
52 and Luo, 2022; Turk-Kubo et al., 2022). The *nifH* gene encodes the identical subunits of the Fe protein of nitrogenase, the  
53 enzyme that catalyzes the  $N_2$  fixation reaction, and contains both highly conserved and variable regions enabling its use as a  
54 phylogenetic marker and as a proxy for  $N_2$ -fixing potential in marine ecosystems globally (Gaby and Buckley, 2011).

55  
56 Although the importance of marine  $N_2$  fixation is well-established, knowledge gaps remain, and discoveries continue to be  
57 made (Zehr and Capone, 2020). For example, high-throughput sequencing (HTS) of *nifH* amplicons is expanding our  
58 knowledge of diazotroph biogeography and activity and has revealed surprising new diversity. However, HTS studies often  
59 utilize different or custom software pipelines and parameters, rendering direct comparisons between studies difficult.  
60 Additionally, many studies do not address the full breadth of diazotrophic diversity because they focus on cyanobacterial  
61 diazotrophs while providing only a superficial analysis of the NCDs present. The resulting lack of information on NCD *in*

62 *situ* distributions limits our understanding of diazotroph ecology and N<sub>2</sub> fixation as well as our ability to predict how these  
63 populations will respond, e.g., trait-based ecological models, to a continually changing ocean.

64

65 To address these issues, we compiled published *nifH* amplicon HTS datasets along with two new datasets. Twenty-one  
66 studies were reprocessed by our newly developed software workflow, which streamlines the integration of multiple, large  
67 amplicon datasets for reproducible analyses. The workflow identifies amplicon sequence variants (ASVs) using a pipeline  
68 developed around DADA2 (Callahan et al., 2016) — the DADA2 *nifH* pipeline — and then executes rigorous post-pipeline  
69 stages to: remove spurious *nifH* ASVs; annotate the remaining quality-filtered ASVs using multiple reference databases and  
70 classification approaches; and obtain *in situ* and modeled environmental data for each sample from the Simons Collaborative  
71 Marine Atlas Project (CMAP; <https://simonscmap.com>). Although created to support research into N<sub>2</sub> fixation (*nifH*), the  
72 complete workflow (ASV pipeline followed by the post-pipeline stages) can be adapted for use with other amplicon datasets,  
73 including other functional genes or taxonomic markers (16S rRNA genes), with some simple modifications.

74

75 In addition to the workflow, our efforts resulted in the construction of a comprehensive database of *nifH* ASVs with  
76 contextual metadata that will be a community resource for marine diazotroph investigations, enhancing comparability  
77 between previous and future *nifH* amplicon datasets. The *nifH* ASV database is available in Figshare  
78 (<https://doi.org/10.6084/m9.figshare.23795943.v2>; Morando et al., 2024a). The entire workflow required to produce the *nifH*  
79 ASV database is available in two GitHub repositories, the DADA2 *nifH* pipeline  
80 ([https://github.com/jdmagasin/nifH\\_amplicons\\_DADA2](https://github.com/jdmagasin/nifH_amplicons_DADA2); Morando et al., 2024b), and the post-pipeline stages  
81 (<https://github.com/jdmagasin/nifH-ASV-workflow>; Morando et al., 2024c).

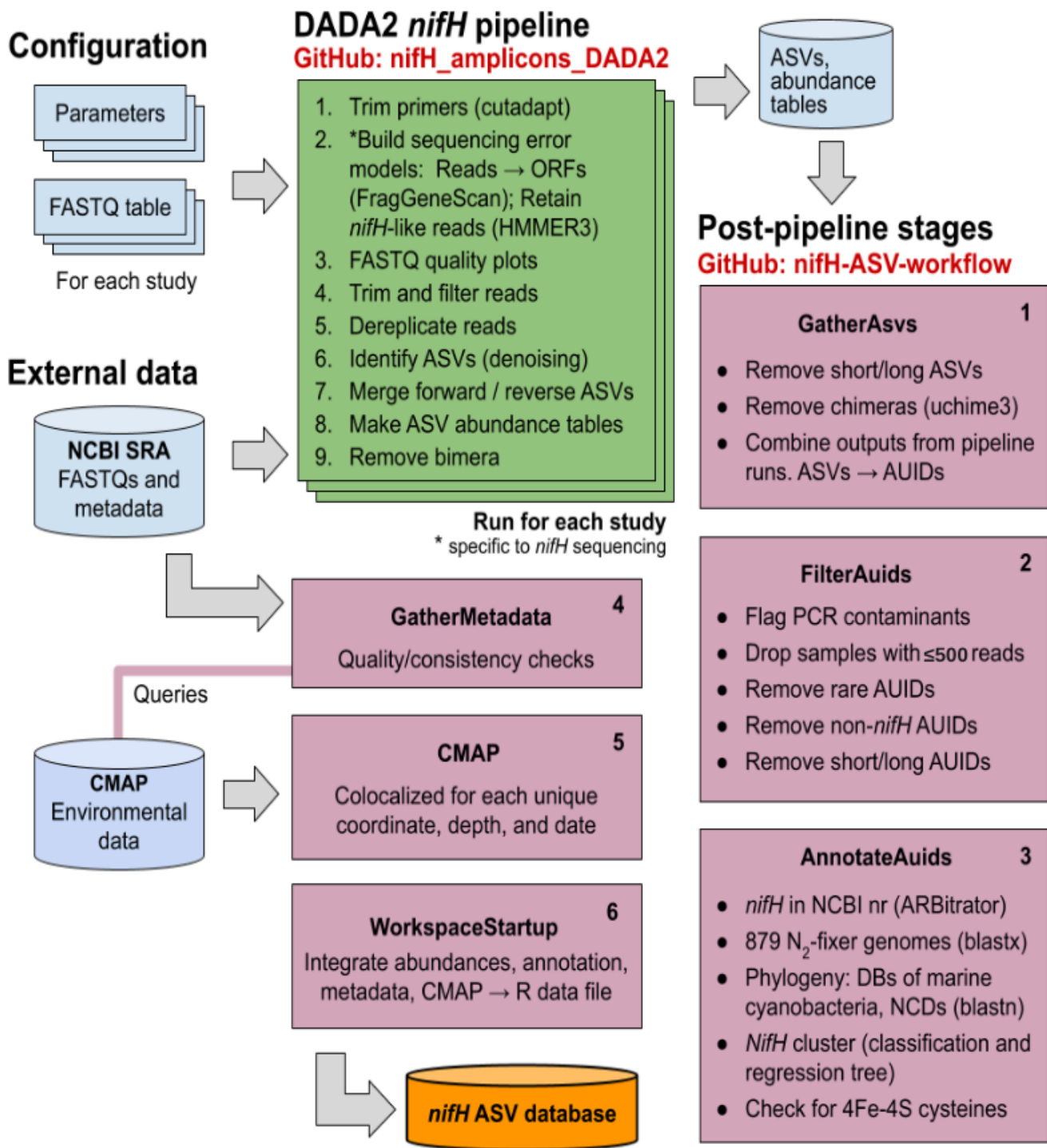
## 82 2 Data and Methods

### 83 2.1 Overview of *nifH* amplicon workflow and *nifH* ASV database generation

84 The full workflow is comprised of two parts: 1) the DADA2 *nifH* pipeline; and 2) a series of post-pipeline stages (Fig. 1).

85

86



87

88 **Figure 1: Schematic of the *nifH* amplicon data workflow.** Data from all studies that met our criteria (Sect. 2.2) were downloaded from  
89 the NCBI Sequence Read Archive (SRA) and processed separately through the DADA2 *nifH* pipeline (green; Sect. 2.3.2), generally using  
90 identical parameters. ASV sequences and abundance tables from all studies were then combined and processed through each stage of the

91 post-pipeline workflow (purple, Sect. 2.3.3) by executing the Makefile associated with each stage. Post-pipeline stages quality-filtered and  
92 then annotated the ASVs by reference to several *nifH* databases (DBs), and downloaded CMAP environmental data matched to the date,  
93 coordinates, and depth of each amplicon dataset. The main output of the entire workflow (pipeline and post-pipeline) is the *nifH* ASV  
94 database, which is available in Figshare (<https://doi.org/10.6084/m9.figshare.23795943.v2>; Morando et al., 2024a). The workflow is  
95 maintained in two GitHub repositories, one for the DADA2 *nifH* pipeline ([https://github.com/jdmagasin/nifH\\_amplicons\\_DADA2](https://github.com/jdmagasin/nifH_amplicons_DADA2);  
96 Morando et al., 2024b) and one for the post-pipeline stages (<https://github.com/jdmagasin/nifH-ASV-workflow>; Morando et al., 2024c).

97

98

99 Required inputs for the pipeline are raw *nifH* amplicon sequencing reads and sample collection metadata (at minimum the  
100 latitude and longitude, depth and sample collection date and time) used to acquire environmental metadata from CMAP.  
101 Criteria for including publicly available datasets are detailed in Section 2.2.1.

102

103 The DADA2 software package is frequently used for processing 16/18S rRNA gene amplicon sequencing data due to its  
104 ability to remove base calling errors (“denoising”) and thereby infer error-free ASVs (Callahan et al., 2016). We have  
105 developed a customizable pipeline to improve the error models utilized by DADA2 by training them only on reads in a  
106 dataset that are valid *nifH* sequences (not PCR artifacts). The DADA2 pipeline runs from the command line in a Unix-like  
107 shell, moving through nine steps (Fig. 1 DADA2 *nifH* pipeline) described in Section 2.3.2 for each study independently.  
108 After the DADA2 pipeline is completed, outputs from all studies are integrated and refined by the six post-pipeline stages of  
109 the workflow, which perform additional quality filtering (e.g., size- and abundance-based selection), identify and remove  
110 spurious sequences (e.g., potential contaminants and non-target sequences), and annotate the ASVs (Fig. 1 Post-pipeline  
111 stages). By considering ASVs from all studies simultaneously, the workflow considers rare ASVs that might be discarded as  
112 irrelevant in a single-study analysis. Workflow stages are executed manually by running their associated Makefiles and  
113 Snakefiles within a Unix-like shell.

114

115 The workflow generates the final data product published in this work, the *nifH* ASV database, which includes ASV  
116 sequences, abundance and annotation tables, sample collection metadata, and sample environmental data from CMAP (Fig.  
117 1). The database is available in Figshare (<https://doi.org/10.6084/m9.figshare.23795943.v2>; Morando et al., 2024a) as a set  
118 of tables (comma-separated value files) and an ASV FASTA file. However, these are also provided within an R data file,  
119 workspace.RData, in the WorkspaceStartup directory in the workflow GitHub repository, for users who wish to analyze,  
120 curate, or customize the database using R packages for ecological analysis. All documentation, scripts, and data needed to  
121 run the workflow and produce the *nifH* ASV database are provided in the workflow GitHub repository  
122 (<https://github.com/jdmagasin/nifH-ASV-workflow>; Morando et al., 2024c). This includes pre-generated pipeline results for  
123 each of the 21 studies as well as the pipeline parameters files.

124

125 In summary, the workflow facilitates the systematic and reproducible exploration of *nifH*-based diversity within microbial  
126 communities and was applied to available *nifH* amplicon data to generate a globally distributed *nifH* ASV database. Together

127 the workflow and *nifH* ASV database will serve as valuable community resources, fostering future investigations while  
128 ensuring comparability between previous and forthcoming studies. In the following sections, detailed descriptions of each  
129 stage of the workflow are provided.

130

## 131 2.2 Compilation of *nifH* amplicon studies

### 132 2.2.1 Published studies

133 We compiled all publicly available *nifH* amplicon HTS data that were generated using the *nifH*1-4 primers (Zani, 1999; Zehr  
134 and McReynolds, 1989) and subsequently sequenced on the Illumina MiSeq/HiSeq platform totaling 21 studies (Table 1).  
135 Limiting the scope to investigations that used the same amplification primers enabled a more tractable comparison across  
136 studies by different research groups that employed varying approaches to sample collection and preparation for sequencing  
137 by different centers. Datasets were downloaded from the National Center for Biotechnology Information (NCBI) Sequencing  
138 Read Archive (SRA) using the GrabSeqs tool (Taylor et al., 2020) by specifying the study's NCBI project accession. Each  
139 dataset obtained included paired-end sequencing reads (in FASTQ files) and a table with the collection metadata for each  
140 sample. Some datasets could not be retrieved directly from the SRA and were obtained from the authors (Table A1). Note  
141 that we did not include studies where data was generated from experimental perturbations or particle enrichments (Table  
142 A1). Data were last accessed from NCBI SRA on 17 April 2024.

143

144 **Table 1: Information on the studies compiled to generate the *nifH* ASV database.** All compiled studies and associated information.  
145 This includes the study ID used to refer to each dataset, the number of samples, NCBI BioProject accession, a reference to each  
146 publication and its corresponding DOI.

147

Study ID	Samples	NCBI BioProject	Reference	DOI
AK2HI	43	PRJNA1062410	This study	n/a
BentzonTilia_2015	56	PRJNA239310	Bentzon-Tilia et al., 2015	10.1038/ismej.2014.119
Ding_2021	32	SUB7406573	Ding et al., 2021	10.3390/biology10060555
Gradoville_2020_G1	111	PRJNA530276	Gradoville et al., 2020	10.1002/lno.11423
Gradoville_2020_G2	56	PRJNA530276	Gradoville et al., 2020	10.1002/lno.11423
Hallstrom_2021	82	PRJNA656687	Hallstrøm et al., 2022b	10.1002/lno.11997
Hallstrom_2022	83	PRJNA756869	Hallstrøm et al., 2022a	10.1007/s10533-022-00940-w
Harding_2018	91	PRJNA476143	Harding et al., 2018	10.1073/pnas.1813658115
Mulholland_2018	29	PRJNA841982	Mulholland et al., 2019	10.1029/2018GB006130
NEMO	56	PRJNA1062391	This study	n/a
Raes_2020	121	PRJNA385736	Raes et al., 2020	10.3389/fmars.2020.00389
Sato_2021	28	PRJDB10819	Sato et al., 2021	10.1029/2020JC017071

<b>Selden_2021</b>	10	PRJNA683637	Selden et al., 2021	10.1002/lno.11727
<b>Shiozaki_2017</b>	22	PRJDB5199	Shiozaki et al., 2017	10.1002/2017GB005681
<b>Shiozaki_2018GBC</b>	20	PRJDB6603	Shiozaki et al., 2018b	10.1029/2017GB005869
<b>Shiozaki_2018LNO</b>	20	PRJDB5679	Shiozaki et al., 2018a	10.1002/lno.10933
<b>Shiozaki_2020</b>	14	PRJDB9222	Shiozaki et al., 2020	10.1038/s41561-020-00651-7
<b>Tang_2020</b>	6	PRJNA554315	Tang et al., 2020	10.1038/s41396-020-0703-6
<b>TurkKubo_2021</b>	136	PRJNA695866	Turk-Kubo et al., 2021	10.1038/s43705-021-00039-7
<b>Wu_2019</b>	18	PRJNA438304	Wu et al., 2019	10.1007/s00248-019-01355-1
<b>Wu_2021</b>	14	PRJNA637983	Wu et al., 2021	10.1007/s10021-021-00702-z

148

149

150

151 Sample quality was validated prior to processing through the DADA2 *nifH* pipeline. Samples were discarded if they did not  
152 contain unmerged pairs of forward and reverse reads with properly oriented primer sequences (Table A1). There were two  
153 exceptions, studies by Shiozaki et al. (2017) and Shiozaki et al. (2018b), that used mixed-orientation sequence libraries and  
154 required preprocessing. The reads in each of these studies were partitioned by whether they captured the coding or template  
155 strand of *nifH*, determined by primer orientation. Because HTS sequence quality generally degrades from 5' to 3', the  
156 partitioned data were run separately through the pipeline to preserve their sequencing error profiles for DADA2. The ASVs  
157 from the misoriented reads (e.g. forward reads with template sequence) were then reverse-complemented and combined with  
158 the properly oriented ASVs into a single ASV abundance table and FASTA file. Table 1 and Table A1 provide information  
159 for obtaining the raw FASTQ files for all samples evaluated for the *nifH* ASV database including information regarding  
160 studies excluded from the database.

161

### 162 2.2.2 Unpublished *nifH* amplicon datasets

163 Additional *nifH* gene HTS datasets were included from DNA samples collected on two cruises in the North Pacific. One was  
164 a transect cruise across the Eastern North Pacific (NEMO; R/V New Horizon, August 2014; Shilova et al., 2017), and the  
165 other was a transect cruise from Alaska to Hawaii (AK2HI; R/V Kilo Moana, September 2017). Euphotic zone samples were  
166 collected from Niskin bottles deployed on a CTD-rosette (NEMO) or from the underway water system (5 m; AK2HI).  
167 NEMO samples (2-4 L) were filtered through 0.2 µm and 3 µm pore-size filters (in series), while AK2HI samples (ca. 2 L)  
168 were filtered through 0.2 µm pore-size filters using gentle peristaltic pumping. Filters were dried, flash frozen and stored at  
169 -80°C until processing. DNA was extracted using a modified DNeasy Plant Kit (Qiagen, Germantown, MD) protocol,  
170 described in detail in Moisander et al. (2008), with on-column washing steps automated by a QIAcube (Qiagen).

171



172 Partial *nifH* DNA sequences were PCR-amplified using the *nifH*1-4 primers in a nested *nifH* PCR assay (Zani, 1999; Zehr  
173 and McReynolds, 1989) according to details in Cabello et al. (2020). All samples were amplified in duplicate and pooled  
174 prior to sequencing. A targeted amplicon sequencing approach was used to create barcoded libraries as described in Green et  
175 al. (2015), using 5' common sequence linkers (Moonsamy et al., 2013) on second round primers, *nifH*1 and *nifH*2. Sequence  
176 libraries were prepared at the DNA Service Facility at the University of Illinois at Chicago, and multiplexed amplicons were  
177 bidirectionally sequenced (2 × 300 bp) using the Illumina MiSeq platform at the W.M. Keck Center for Comparative and  
178 Functional Genomics at the University of Illinois at Urbana-Champaign. Samples were multiplexed to achieve ca. 40,000  
179 high quality paired reads per sample. The AK2HI and NEMO datasets can be found in the SRA (BioProjects  
180 PRJNA1062410 and PRJNA1062391, respectively).

181

### 182 2.2.3 Sample collection data and co-localized CMAP environmental data

183 Sample collection data (e.g. coordinates, depth, date) and environmental data provide essential context for the interpretation  
184 of diazotroph 'omics datasets. Large-scale multivariate analyses depend on properly formatted, complete, and ideally quality  
185 checked metadata from consistently collected and analyzed measurements. However, accessibility to this information is often  
186 limited (especially environmental data) for datasets published across multiple decades. Therefore, we first obtained sample  
187 collection metadata from the SRA, and corrected or flagged errors and inconsistencies in the GatherMetadata stage of our  
188 post-pipeline workflow (described below), to ensure consistency and completeness. For each sample, the geographic  
189 coordinates, depth, and collection date (at local noon) from the SRA were used to query the Simons Collaborative Marine  
190 Atlas Project on 24 March 2023 (CMAP; <https://simonscmap.com/>; Ashkezari et al., 2021) for co-localized environmental  
191 data using a custom script (`query_CMAP.py`) in the CMAP stage of the workflow (Fig. 1). CMAP is an open-source data  
192 portal designed for retrieving, visualizing, and analyzing diverse ocean datasets including research cruise-based and  
193 autonomous measurements of biological, chemical, and physical properties, multi-decadal global satellite products, and  
194 output from global-scale biogeochemical models. For each sample a mixture of 100 satellite derived and modeled  
195 environmental variables from the CMAP repository were obtained. These, along with the SRA collection data, are included  
196 in our database. Aggregated metadata for all samples are summarized in Supplementary Table 1 but a detailed description of  
197 environmental metadata can be found at the CMAP website (<https://simonscmap.com/catalog>). Metadata are available in the  
198 *nifH* ASV database (`metaTab.csv` for sample metadata and `cmapTab.csv` for environmental data).

199



## 200 2.3 Automated workflow for processing datasets with the DADA2 *nifH* pipeline

### 201 2.3.1 Installation of the DADA2 *nifH* pipeline and the post-pipeline workflow

202 The workflow (Fig. 1) comprises two software projects installed from separate GitHub repositories,  
203 *nifH\_amplicons\_DADA2* which contains the ASV pipeline and ancillary scripts, and *nifH-ASV-workflow* which integrates  
204 pipeline results for all datasets with annotation and CMAP environmental data to produce the data deliverable of the present  
205 work, the *nifH* ASV database. Installation requires cloning the *nifH\_amplicons\_DADA2* repository  
206 ([https://github.com/jdmagasin/nifH\\_amplicons\\_DADA2](https://github.com/jdmagasin/nifH_amplicons_DADA2); Morando et al., 2024b) to a local machine and then downloading  
207 several external software packages using miniconda3. Detailed installation instructions are available from the GitHub  
208 homepage, as well as a small tutorial to verify the installation on a small *nifH* amplicon dataset and introduce the two main  
209 pipeline commands (`organizeFastqs.R` and `run_DADA2_pipeline.sh`). Altogether the installation and example take 30–40  
210 min.

211

212 After installing the ASV pipeline, installation of the *nifH-ASV-workflow* proceeds similarly: Clone the GitHub repository  
213 (<https://github.com/jdmagasin/nifH-ASV-workflow>; Morando et al., 2024c) and then download a few additional packages  
214 with miniconda3 (~10 min to complete). For each study, the *nifH-ASV-workflow* includes the pipeline outputs (ASVs and  
215 abundance tables) which were used to create the *nifH* ASV database. Pipeline parameters and FASTQ input tables for each  
216 study are also provided for users who instead wish to rerun the pipeline starting from FASTQs downloaded from the SRA.  
217 Because the *nifH-ASV-workflow* includes data and parameters specific to the studies used in this work, it has a separate  
218 GitHub repository from the pipeline. However, we emphasize that together they comprise the *nifH* amplicon workflow in  
219 Fig. 1.

220

221 Adding a new dataset to the workflow can be summarized in four steps: (1) Start a Unix-like shell that includes the required  
222 software (by “activating” a miniconda3 environment called `nifH_ASV_workflow`). (2) Generate ASVs for the new dataset by  
223 running it through the pipeline, likely multiple times to tune parameters (Table 2). Output can be placed in the Data directory  
224 alongside other studies used in this work, and SRA metadata must be added to `Data/StudyMetadata`. (3) Include the new  
225 ASVs in the workflow by appending rows to the table `GatherASVs/asvs.noChimera.fasta_table.tsv`, which has file paths to  
226 all ASV abundance tables. (4) For each stage shown in Fig. 1, run the associated Makefile or Snakefile from the Unix-like  
227 shell by executing `"make"` or `"snakemake -c1 --use-conda"`, respectively. Documentation resides within each Makefile or  
228 Snakefile. Input tables from the post-pipeline workflow also have embedded documentation.

229

230 **Table 2. Parameters for controlling the DADA2 *nifH* pipeline.** Default values can be overridden in the text file that is passed to  
231 `run_DADA2_pipeline.sh`. Parameters for "Read trimming" and "Error models" are used in steps 1 and 2 of the pipeline (Fig. 1). The  
232 remaining parameters are important for controlling how DADA2 trims and quality filters the reads, and merges forward and reverse  
233 sequences to create ASVs.

DADA2 <i>nifH</i> pipeline step	Parameter name	Default value	Description	Studies with non-default parameters
<b>Read Trimming</b> Remove primers with cutadapt	forward	TGYGAYCCN AARGCNGA	Forward primer 5' to 3'. Default is <i>nifH2</i> (Zehr and McReynolds, 1989).	None
	reverse	ADNGCCATC ATYTCNCC	Reverse primer 5' to 3'. Default is <i>nifH1</i> (Zehr and McReynolds, 1989).	None
	allowMissingPrimers	FALSE	If TRUE, retain read pairs even if primers are absent, e.g. if trimmed reads were uploaded to NCBI SRA.	Ding et al., 2021
<b>Error Models</b>	skipNifHErrorModels	FALSE	By default, use only <i>nifH</i> -like reads to train error models. If TRUE, use a random sample of all reads.	None
	NifH_minBits	150	Train error models using reads that align to PFAM00142 at $\geq$ the specified bit score. The trusted cut off in PFAM00142 (25 bits) is always used to filter reads, then NifH_minBits. If set to 0, only the trusted cut off is used.	Set to 0 for most studies. Exceptions that used 100 bits were: Bentzon-Tilia et al., 2015; Gradoville et al., 2020; Shiozaki et al., 2018a; Turk-Kubo et al., 2021.
	NifH_minLen	33	Train error models using reads with ORFs that align with $\geq$ this many residues to PFAM00142.	None
<b>DADA2 filterAndTrim</b>	id.field	NA	Specify number of ID field if reads do not follow the CASAVA format. Forwarded to filterAndTrim(). If set, usually to 1.	Ding et al., 2021; Wu et al., 2021; Wu et al., 2019; Mulholland et al., 2019; Raes et al., 2020; Tang et al., 2020; Selden et al., 2021; Hallstrøm et al., 2022b; Hallstrøm et al., 2022a
	maxEE.fwd	Inf	Forwarded to filterAndTrim().	All studies set to 2.
	maxEE.rev	Inf		All studies set to 4.
	minLen	20	Forwarded to filterAndTrim().	None
	truncLen.fwd	0	Forwarded to filterAndTrim().	Ancillary script estimateTrimLengths.R determined optimal lengths.
	truncLen.rev	0		
	truncQ	2	Forwarded to filterAndTrim()	All studies used truncLen.
useOnlyR1Reads	FALSE	If TRUE, only use R1 reads (and do not call mergePairs()). Used if R2 reads are very low quality.	None	
<b>DADA2 mergePairs</b>	minOverlap	12	Forwarded to mergePairs().	None
	maxMismatch	0	Forwarded to mergePairs().	All studies set to 1.
	justConcatenate	FALSE	Forwarded to mergePairs().	None

234

235

### 236 2.3.2 DADA2 *nifH* pipeline

237 To encourage reproducible outputs and usage by non-programmers, the DADA2 pipeline (GitHub repository:  
238 *nifH\_amplicons\_DADA2*) is controlled by a plain text parameters file (Table 2) and a descriptive table of input samples (the  
239 “FASTQ map”). Since a study might include samples with vastly different diazotroph communities and relative abundances,

240 potentially impacting ASV inferences by DADA2, the FASTQ map for a study enables samples to be partitioned into  
241 "processing groups" that are each run separately through DADA2. For example, in the present work processing groups  
242 usually partitioned the samples in a study by the unique combinations of collection station or date, nucleic acid type (DNA or  
243 RNA), size fraction, and collection depth. Pipeline outputs for each processing group are stored in a directory hierarchy with  
244 levels that follow the processing group definition. Partitioning datasets into processing groups greatly improves the overall  
245 speed of DADA2 and simplifies subsequent analyses that compare ASVs detected in different kinds of samples (e.g.,  
246 detected versus transcriptionally active diazotrophs, or presence across different stations, depths, and/or size fractions). For  
247 generating the *nifH* ASV database, studies that met selection criteria (Sect. 2.2.1 and Table 1) were run through the pipeline  
248 using the study-specific FASTQ maps and parameters available in the Data directory of the *nifH*-ASV-workflow GitHub  
249 repository.

250

251 The DADA2 pipeline runs from the command line in a Unix-like shell, moving through 9 main steps (Fig. 1 DADA2 *nifH*  
252 pipeline): (1) trim reads of primers using cutadapt (Martin, 2011); (2) build sequencing error models; (3) make FASTQ  
253 quality plots; (4) trim and filter reads based on quality; (5) dereplicate; (6) denoise (ASV inference); (7) merge forward and  
254 reverse sequences; (8) make the ASV abundance table; and (9) remove bimeras (Callahan et al., 2016 for steps 2 through 9).  
255 These steps will be familiar to DADA2 users, except that for step 2 the error models are trained only on *nifH*-like reads  
256 (discussed below). To run the pipeline on other functional genes, the parameters file would need to be edited to disable  
257 *nifH*-based error models and to include the expected primers. We again note that the DADA2 pipeline is distinct from the  
258 post-pipeline workflow stages which are specific to this work, but together they comprise the workflow in Fig. 1.

259

260 DADA2 parameters impact the ASV sequences identified, and the number of reads used. Thus, exploring parameters is  
261 essential for checking the robustness of ASVs (particularly rare ones) and their relative abundances. The method and  
262 parameters used to trim the reads are especially important because most pipeline steps occur after filterAndTrim (Fig. 1).  
263 Two methods are supported: One can trim each read based on its quality degradation (truncQ parameter to the DADA2  
264 filterAndTrim function; Table 2) or all reads at the same position determined by inspecting sample FASTQ quality plots  
265 (truncLen parameter; Table 2, and comparison of methods in Appendix B). The latter approach can be labor-intensive and  
266 unsystematic for studies with tens to hundreds of samples. To address this the ancillary script estimateTrimLengths.R can be  
267 used to determine trimming lengths that will maximize the percentage of reads that make it through the pipeline. For each  
268 FASTQ file in a study, the script chooses 1 K read pairs at random and removes the primers. Then the read pairs are trimmed  
269 using every combination of lengths over a window (from 55—85 % of the median read length in 15 bp steps) and successful  
270 merges (with  $\geq 12$  bp overlapping and  $\leq 2$  mismatches) are counted. The counts are averaged across all samples (weighting by  
271 sequencing depths) and the top ten combinations of forward and reverse trimming lengths are reported in a table, with  
272 estimates for the percentages of reads retained and the mean errors per read to help choose the maxEE parameters (Table 2).

273

274 The pipeline allows one to rerun DADA2 steps 3–9, with outputs saved in separate, date-stamped directories. Primer  
275 removal and error models (steps 1–2) are unlikely to benefit much from parameter tuning, so the pipeline reuses outputs  
276 from those steps. Log files and diagnostic plots created by the pipeline are intended to facilitate parameter evaluation as well  
277 to capture statistics to support publication. Moreover, logs and other pipeline outputs are consistently formatted across  
278 pipeline runs, which enables scripts to aggregate and analyze results across datasets such as in our workflow.

279

280 Step 1 consisted only of primer removal using cutadapt (Martin, 2011). Raw reads were retained only if the forward (nifH2)  
281 and reverse (nifH1) primers were both found on the R1 and R2 reads, respectively. DADA2 sequencing error models were  
282 built at step 2 using only the reads predicted to be *nifH*, rather than a subsample of all reads as in typical use of DADA2.  
283 Reads likely to encode *nifH* were identified as follows: FragGeneScan (version 1.31, (Rho et al., 2010)) was used to predict  
284 open reading frames (ORFs) on R1 reads which were then aligned to the nitrogenase PFAM model (PF00142.20) using  
285 HMMer3 (hmmsearch version 3.3.2; hmmer.org). ORFs with >33 residues and a bit score that exceeded the trusted cut-off  
286 encoded in the model (25.0 bits) were retained. Prefiltering the reads aims to reduce effects of PCR artifacts on the error  
287 models. For some studies this approach resulted in increases (~3–10 %) in the total percentage of reads retained in ASVs,  
288 and fewer total ASVs, compared to using error models based on a subsample of all reads. Adapting the pipeline to a different  
289 marker gene would only require substituting an appropriate PFAM model, or disabling step 2 (by setting  
290 skipNifHErrorModels to TRUE; Table 2), which forces the pipeline to make error models by subsampling from all reads. At  
291 step 4, DADA2 filterAndTrim() trimmed forward and reverse reads using the lengths suggested by estimateTrimLengths.R  
292 and then discarded read pairs that had excessive errors (>2 for R1 reads, >4 for R2 reads) or were <20 bp. Conservative  
293 parameters were used for merging sequences: At most 1 base pair was allowed to mismatch in the forward and reverse  
294 sequence overlap of minimally 12 bp (stage 7). Dereplicating (step 5) and denoising, ASV calling (step 6), generating an  
295 abundance table (step 8), and bimera detection (step 9), were all performed with default DADA2 parameters. Datasets that  
296 passed pre-processing steps (Table 1) were run through the DADA2 pipeline using mostly identical parameters except for the  
297 trimming lengths (truncLen.fwd and truncLen.rev in Table 2).

298

### 299 2.3.3 Post-pipeline stages

300 The workflow post-pipeline stages (GitHub repository: nifH-ASV-workflow) combine the pipeline outputs, conduct further  
301 quality control steps, co-locate the samples with environmental data from the CMAP data portal, and annotate the ASVs  
302 (Fig. 1 Post-pipeline stages). Key outputs from the post-pipeline are: a unified FASTA with all the unique ASVs detected  
303 across all the studies (i.e. all samples); tables of ASV total counts and relative abundances in all studies; multiple annotations  
304 for each ASV by comparison to several *nifH* reference databases; and CMAP environmental data for each sample. These  
305 outputs comprise the *nifH* ASV database, and are all available within an R image file (workspace.RData) generated by the

306 workflow which is included in the nifH-ASV-workflow repository. Provision as an R image will make the outputs  
307 immediately accessible to many researchers who prefer R due to its extensive packages for ecological analysis. The *nifH*  
308 ASV database is also available on Figshare (<https://doi.org/10.6084/m9.figshare.23795943.v2>; Morando et al., 2024a). The  
309 remainder of this section describes each of the post-pipeline stages.

310

311 The GatherAsvs stage aggregated ASV sequences and abundances across all DADA2 pipeline runs (i.e. from all samples and  
312 studies). First, ASVs were filtered based on length. Chimera sequences were then removed using UCHIME3 denovo (Edgar,  
313 2016a) via VSEARCH (Rognes et al., 2016). Chimera sequences were identified within each sample, but the final  
314 classification was based on majority vote (chimera or not) across the samples in the processing group. Second, the  
315 GatherAsvs stage combined the non-chimeric ASVs from all studies into a single abundance table and FASTA file. Since  
316 each study is run independently through the DADA2 pipeline, ASV identifiers are not consistent across studies. Therefore,  
317 each unique ASV sequence was renamed with a new unique identifier of the form AUID.*i*, where AUID stands for **ASV**  
318 **U**niversal **I**dentifier. The scripts used to rename the ASVs (assignAUIDs2ASVs.R) and to create the new abundance table  
319 (makeAUIDCountTable.R) are available at the nifH\_amplicons\_DADA2 GitHub repository (in  
320 scripts.ancillary/ASVs\_to\_AUIDs). The script assignAUIDs2ASVs.R optionally takes an AUID reference FASTA so that  
321 AUIDs can be preserved as new datasets are added to future versions of the *nifH* ASV database.

322

323 Both rare and potential non-*nifH* sequences were assessed on the unified AUID tables in the next stage, FilterAuids (Fig. 1).  
324 First, possible contaminants were identified by the Makefile invocation of check\_nifH\_contaminants.sh, provided as an  
325 ancillary script in the pipeline GitHub repository. In brief, check\_nifH\_contaminants.sh first translated all ASVs into amino  
326 acid sequences using FragGeneScan (Rho et al., 2010), which were then compared using *blastp* to 26 contaminants known  
327 from previous *nifH* amplicon studies (Zehr et al., 2003; Goto et al., 2005; Farnelid et al., 2009; Turk et al., 2011). ASVs that  
328 aligned at >96 % amino acid identity to known contaminants were flagged. Next FilterAuids removed samples with ≤500  
329 reads, and rare ASVs, defined as those that did not have at least one read in at least two samples or ≥1000 reads in one  
330 sample.

331

332 Next, the ancillary script, classifyNifH.sh, was employed to identify and remove non-*nifH*-like sequences. The script utilized  
333 *blastx* to search each ASV against ~44 K positive and ~15 K negative examples of NifH protein sequences that were found  
334 in NCBI GenBank by ARBitrator (run on April 28, 2020; Heller et al., 2014). ASVs were classified based on the relative  
335 quality of their best hits in the two databases, similar to the "superiority" check in ARBitrator. An ASV was classified as  
336 positive if the E-value of its best positive hit was ≥10 times smaller than the E-value for the best negative hit, and vice versa  
337 for negative classifications. ASVs failing to meet these criteria were classified as 'uncertain'. The *blastx* searches used the  
338 same effective sizes for the two databases (-dbsize 1000000), so that E-values could be compared, and retained up to 10 hits  
339 (-max\_target\_seqs 10).

340

341 The FilterAuids stage of the workflow exclusively discarded ASVs with negative classifications. “Uncertain” ASVs were  
342 retained as potential *nifH* sequences not in GenBank. In the last stage, FilterAuids excluded ASVs with lengths that fell  
343 outside 281–359 nucleotides, a size range which in our experience encompasses the majority of valid *nifH* amplicon  
344 sequences generated by nested PCR with nifH1–4 primers.

345

346 For each AUID in the *nifH* ASV database, we provide taxonomical annotations using several different approaches,  
347 encompassed by the AnnotateAuids stage (Fig. 1) and accessible through ancillary scripts in the GitHub repository (in  
348 scripts.ancillary/Annotation). The script blastxGenome879.sh enables a protein level comparison via *blastx* against a  
349 database of 879 sequenced diazotroph genomes (“genome879”, <https://www.jzehrlab.com/nifH>). Here, the closest cultivated  
350 relative for each AUID was determined by smallest E-value among alignments with  $\geq 50$  % amino acid identity and  $\geq 90$  %  
351 query sequence coverage. Cautious interpretation is suggested because the reference database is small and contains only  
352 cultivable taxa. Similarly, the top nucleotide match of each AUID was identified by E-value within alignments possessing  
353  $\geq 70$  % nt identity and  $\geq 90$  % query sequence coverage obtained by *blastn* against a curated database of *nifH* sequences (July  
354 2017 *nifH* database, <https://www.jzehrlab.com/nifH>) by executing the blastnARB2017.sh script. Additionally, *nifH* cluster  
355 annotations were assigned to each ASV using the classification and regression tree (CART) method of Frank et al. (2016).  
356 This approach was implemented as part of a custom tool that predicted ORFs for the ASVs with FragGeneScan, then  
357 performed a multiple sequence alignment on the ORFs, and then applied the CART classifier. The tool is available as the  
358 ancillary script assignNifHclustersToNuclSeqs.sh.

359

360 The Makefile created and searched against two “phyloTYPE” databases, one containing 223 *nifH* sequences from prominent  
361 marine diazotrophs including NCDs (Turk-Kubo et al., 2022) and another with 44 UCYN-A *nifH* oligotype sequences  
362 (Turk-Kubo et al., 2017). These databases were searched using *blastn* with the effective database size of the ARB2017  
363 database (-dbsize set to ~29 million bases) to enable E-value comparisons across all three searches. For each ASV, we  
364 provide phyloTYPE annotations based on the top hit by E-value if the alignment had  $\geq 97$  % nt identity and covered  $\geq 70$  % of  
365 the ASV. Finally, ORFs for all ASVs were searched for highly conserved residues which are thought to coordinate the  
366 4Fe-4S cluster in NifH, specifically for paired cysteines shortly followed by AMP residues (described in Schlessman et al.  
367 1998). This simple check, performed by the script check\_CCAMP.R, was intended to complement the reference-based  
368 annotations above. Presence of cysteines and AMP could be used to retain ASVs that have no close reference. Absence could  
369 be used to flag ASVs that, despite high similarity to a reference sequence, might not represent functional *nifH* (e.g. due to  
370 frameshifts).

371

372 Since the annotation scripts provided multiple taxonomic identifications for most of the AUIDs, a primary taxonomic ID was  
373 assigned for each AUID using the script make\_primary\_taxon\_id.py. If a phyloTYPE annotation (e.g., Gamma A) was  
374 assigned, this became the primary taxonomic ID; otherwise, cultivated diazotrophs from genome879 were used (e.g.,

375 “*Pseudomonas stutzeri*”). Finally, when neither a phylotype nor a cultivated diazotroph could be determined, the *nifH* cluster  
376 (e.g. “unknown 1G”) was used. AUIDs without an assigned *nifH* cluster or taxonomic rank below domain were removed  
377 from the final *nifH* ASV database unless paired cysteines and AMP were detected. This final data filtration step occurred in  
378 the WorkspaceStartup stage described below.

379

380 The CMAP stage was managed by a Snakefile that called the script query\_cmap.py to query the CMAP data portal for  
381 co-localized environmental data (Fig. 1). The script was passed the main output from the GatherMetadata stage,  
382 metadata.cmap.tsv, a table of the collection coordinates, dates at local noon, and depths from all the samples.  
383 GatherMetadata reported any samples with missing metadata and ensured standardized formats for the required query fields.  
384 Additionally, query\_cmap.py validated fields prior to querying CMAP. It should be noted that the precision of values  
385 obtained from CMAP depend on floating point arithmetic, not the significant digits of the underlying measurement or model.  
386 Therefore, prior to an analysis requiring high precision for specific CMAP variables, it is recommended to consult the  
387 original producer of the data to determine the significant digits.

388

389 The last stage of the workflow, WorkspaceStartup, filtered out AUIDs that had no annotation and then generated the final  
390 *nifH* ASV database, which is comprised of AUID abundance tables (counts and relative), AUID annotations, sample  
391 metadata and corresponding environmental data. These data are provided as text files (.csv and FASTA) within a single  
392 compressed file (.tgz) that is available in Figshare (<https://doi.org/10.6084/m9.figshare.23795943.v2>; Morando et al., 2024a)  
393 as well as within the workflow GitHub repository within an R image file (workspace.RData).

#### 394 **2.4 Diazotroph biogeography from DNA dataset of the *nifH* ASV database**

395 The DNA dataset, a custom version of the *nifH* ASV database restricted to DNA samples (representing a majority of the  
396 database, only removing 108 cDNA samples out of 944 total samples), was created to showcase the utility of the workflow.  
397 Additional data reduction steps were conducted, averaging replicates and samples from the same location but different size  
398 fractions, to enable comparisons between different sampling methodologies.

### 399 **3 Results and Discussion**

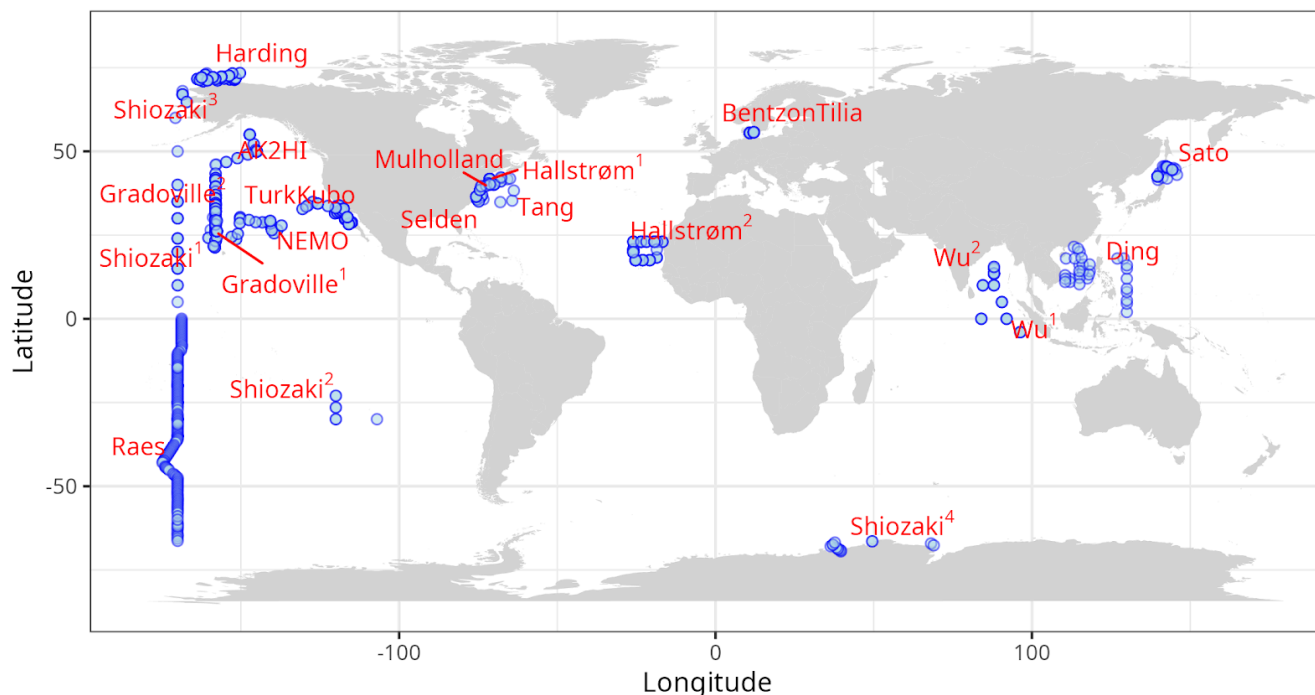
#### 400 **3.1 Generation of the marine *nifH* ASV database**

401 All publicly available marine *nifH* amplicon HTS data from studies that met our criteria, including two new studies, were  
402 compiled in the present investigation (see Sect. 2.2 and Table A1). Altogether 982 samples from 21 studies, comprising a  
403 total of 87.7 million reads (Table 3), were processed through the entire workflow, i.e., the DADA2 *nifH* pipeline (Sect. 2.2.2)  
404 as well as the post-pipeline stages (Sect. 2.2.3). The *nifH* ASV database, i.e., the ASV sequences, abundances, and  
405 annotations, as well as sample collection and CMAP environmental data, was generated from the 944 samples, 9383 ASVs,



406 and 43.0 million reads that were retained by this workflow (Figs. 1 and 2 and Table 3). To our knowledge it is the only global  
 407 database for marine diazotrophs detected using *nifH* HTS amplicon sequencing, with comprehensive, standardized ancillary  
 408 data (Fig. 2 and Supplementary Table 1).

409



410

411 **Figure 2: Global sampling distribution of the *nifH* ASV database.** World map of sampling locations for the datasets compiled and  
 412 processed to construct the *nifH* ASV database. Abbreviated study IDs are used with superscripts ordered by publication year for Shiozaki  
 413 (2017, 2018GBC, 2018LNO, and 2020), Hallström (2021 and 2022), and Wu (2019 and 2021). For Gradoville the superscripts indicate  
 414 Gradients cruises 1 and 2. See Table 1 for the citation source linked to each study ID.

415

416

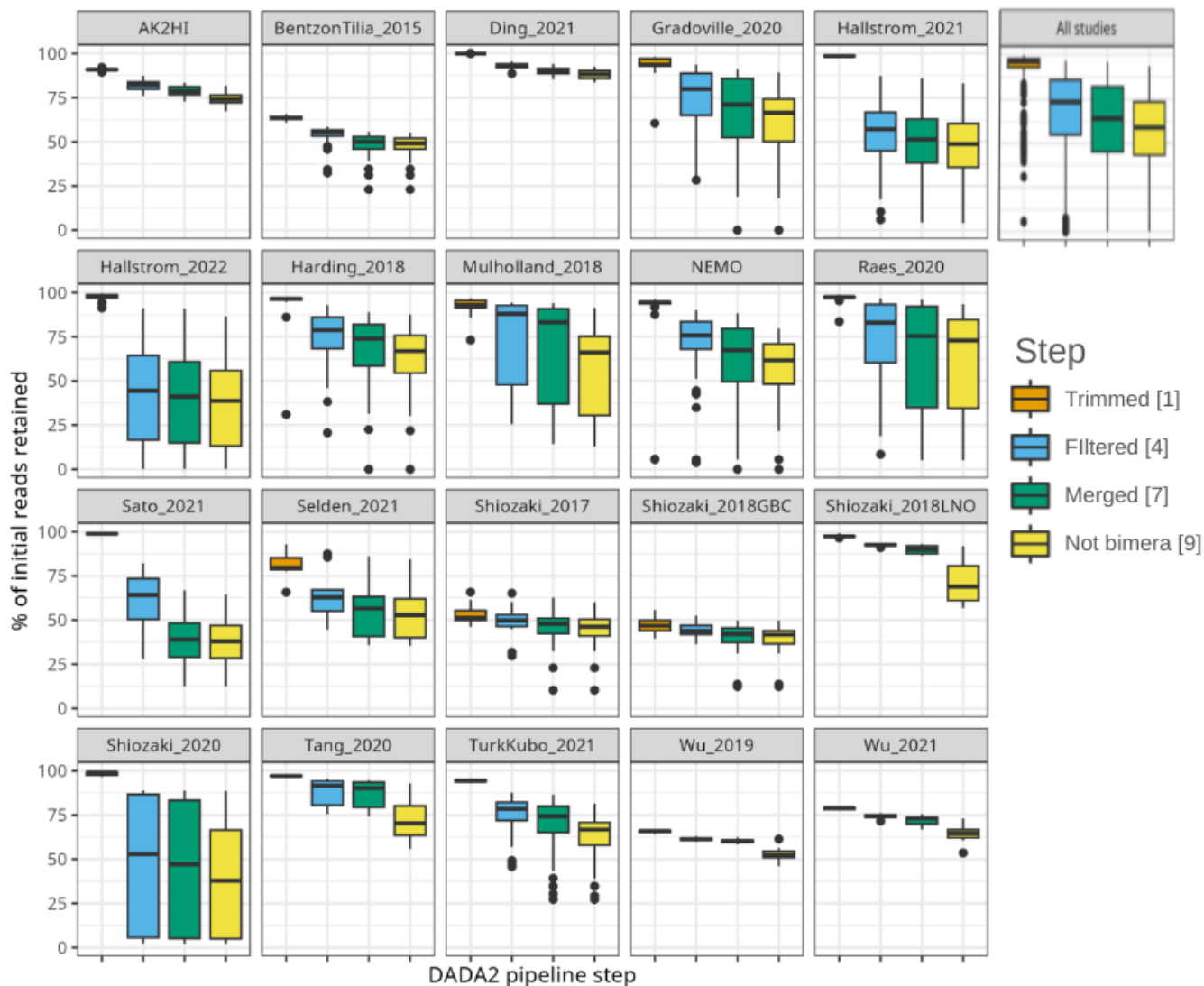
417 **Table 3: Summary of the full *nifH* workflow.** The number of samples, ASVs, and reads retained through the entire workflow (the  
 418 DADA2 *nifH* pipeline and major post-pipeline stages) to create the *nifH* ASV database. The vast majority ASVs that were removed by  
 419 GatherAsvs fell outside 200–450 nt. WorkspaceStartup removed ASVs with no annotation and samples that had zero reads after ASV  
 420 filtering.

	Initial	DADA2 pipeline	Gather Asvs	FilterAuids				Workspace Startup
				≤500 reads in sample	rare	non-NifH	length	
<b>Samples</b>	982	982	982	951	951	951	951	944
<b>ASVs</b>	n/a	152,915	139,355	139,334	18,193	16,253	11,915	9,383
<b>Reads (millions)</b>	87.7	48.7	48.4	48.4	45.5	45.0	43.8	43.0

421

422

423 Interestingly, studies were affected differently by each step of the DADA2 *nifH* pipeline (Fig. 3 and Table 4). There were  
 424 major losses of reads during ASV merging, with several studies retaining <40 % of their total reads by the end of the pipeline  
 425 (i.e., Hallstrom\_2022, Sato\_2021, and Shiozaki\_2020), though on average about 60 % of the reads were retained across  
 426 studies (Fig. 3 and Table 4).  
 427



428

429 **Figure 3: Study-specific retention of reads at each stage of the pipeline.** The proportion of total reads in each sample that are retained  
 430 at the completion of each step of the DADA2 *nifH* pipeline. Each box shows the distribution for samples in the indicated study (using  
 431 Study IDs in Table 1), or for all samples together (top right). Proportions for Shiozaki\_2017 and Shiozaki\_2018GBC reflect that  
 432 approximately half the amplicons were not in the orientation expected by the pipeline (see text). Numbers in the legend indicate pipeline  
 433 steps in Figure 1.

434

435 **Table 4: Quality filtering by the DADA2 *nifH* pipeline.** For each study ID are shown the mean numbers of reads retained per sample at  
 436 the end of each stage of the DADA2 *nifH* pipeline, as well as the mean percentage of reads retained. Statistics in the bottom three rows  
 437 pool all samples. Initial, Trimmed<sup>4</sup>, Filtered<sup>4</sup>, and Merged<sup>9</sup> and non-Bimera<sup>9</sup> and their superscripts are specific to the pipeline steps in  
 438 Figure 1. At each step (column) the calculations include only the samples that have >0 reads.

Study		Initial	Trimmed <sup>4</sup>	Filtered <sup>4</sup>	Merged <sup>9</sup>	Non-bimera <sup>9</sup>	Retained (%)
AK2HI		4.5E+04	4.1E+04	3.7E+04	3.6E+04	3.3E+04	74.1
BentzonTilia_2015		8.2E+03	5.2E+03	4.6E+03	4.1E+03	4.1E+03	48.1
Ding_2021		5.6E+04	5.6E+04	5.2E+04	5.0E+04	4.9E+04	88.1
Gradoville_2020		4.0E+04	3.8E+04	2.9E+04	2.6E+04	2.4E+04	60.3
Hallstrom_2021		2.5E+05	2.5E+05	1.5E+05	1.4E+05	1.4E+05	48.7
Hallstrom_2022		2.0E+05	1.9E+05	7.5E+04	7.4E+04	6.6E+04	36.3
Harding_2018		4.2E+04	4.1E+04	3.1E+04	2.9E+04	2.6E+04	63.2
Mulholland_2018		1.8E+05	1.6E+05	1.3E+05	1.2E+05	1.0E+05	58.5
NEMO		5.7E+04	5.4E+04	4.2E+04	3.6E+04	3.3E+04	57.1
Raes_2020		9.3E+04	9.1E+04	7.7E+04	6.9E+04	6.5E+04	61.0
Sato_2021		7.5E+04	7.4E+04	4.5E+04	2.9E+04	2.9E+04	38.8
Selden_2021		1.5E+05	1.2E+05	9.2E+04	8.2E+04	8.0E+04	54.7
Shiozaki_2017		1.8E+04	9.3E+03	8.9E+03	8.4E+03	8.2E+03	44.1
Shiozaki_2018GBC		2.4E+04	1.1E+04	1.1E+04	1.0E+04	9.8E+03	38.6
Shiozaki_2018LNO		6.7E+04	6.5E+04	6.2E+04	6.0E+04	4.8E+04	71.5
Shiozaki_2020		2.5E+05	2.5E+05	1.8E+05	1.8E+05	1.4E+05	39.1
Tang_2020		4.7E+04	4.6E+04	4.1E+04	4.0E+04	3.4E+04	72.4
TurkKubo_2021		5.5E+04	5.2E+04	4.2E+04	4.0E+04	3.6E+04	63.2
Wu_2019		8.0E+04	5.3E+04	4.9E+04	4.8E+04	4.2E+04	52.9
Wu_2021		8.0E+04	6.3E+04	6.0E+04	5.8E+04	5.2E+04	64.4
All samples and studies	mean	8.9E+04	8.5E+04	5.8E+04	5.4E+04	4.9E+04	56.9
	median	5.1E+04	4.8E+04	3.7E+04	3.2E+04	3.0E+04	59.0
	sum	8.8E+07	8.4E+07	5.7E+07	5.3E+07	4.8E+07	60.0

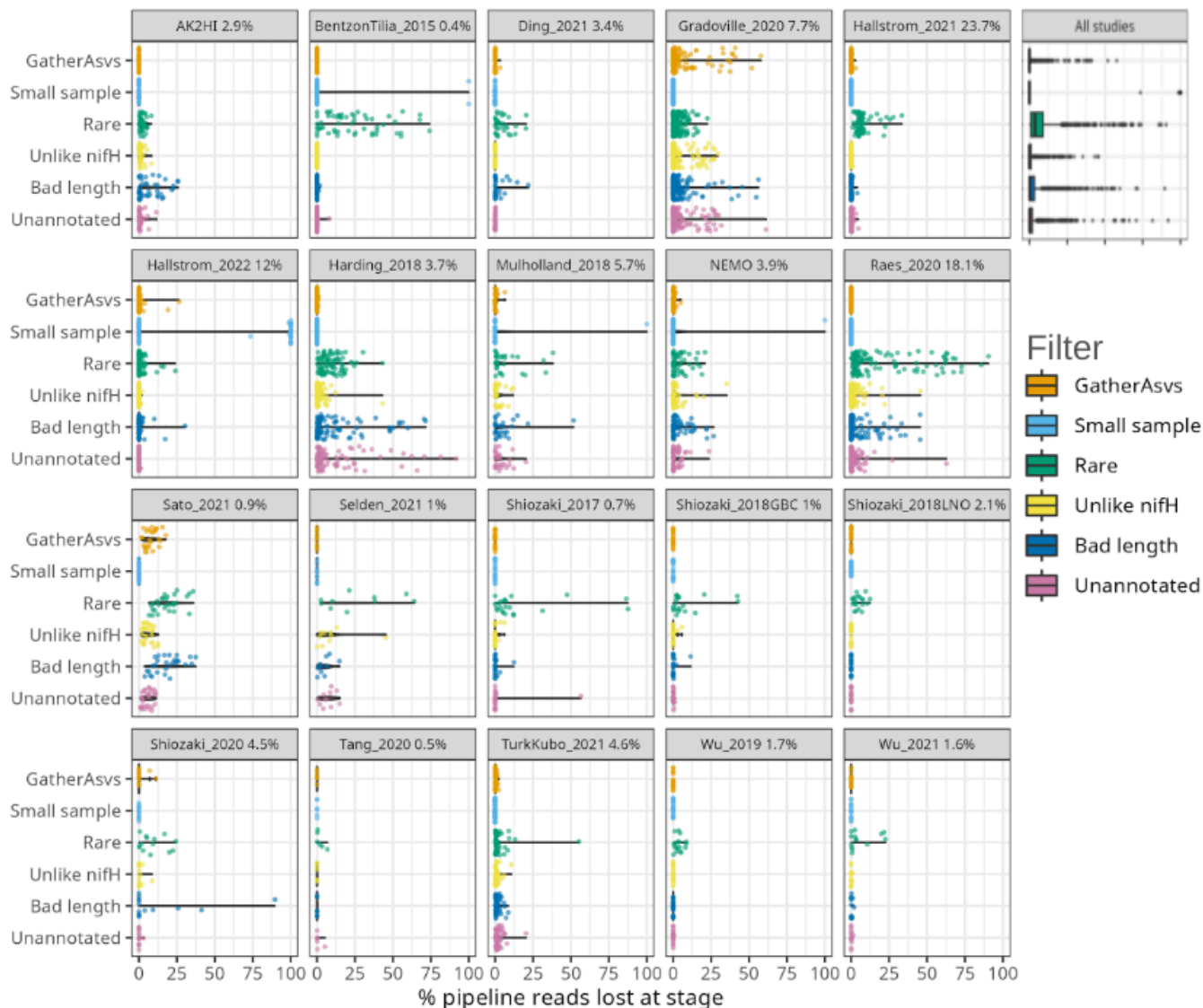
439

440

441

442 Post-pipeline stages of the workflow further refined the data (detailed in Methods) (Fig. 4). First, GatherAsvs identified and  
 443 removed 163 chimeras using uchime3 denovo (distinct from the bimera filtering done by the pipeline), and then removed 8.7

444 K ASVs that were far outside expected *nifH* lengths (200–450 nt). AUIDs were assigned to the remaining 139 K unique  
 445 non-chimeric ASVs (comprising 48.4 million total reads; Tables 3 and 5). The FilterAuids stage had the largest impacts on  
 446 retained data. Thirty-one samples with  $\leq 500$  reads were removed because they would likely misrepresent their diazotrophic  
 447 communities. The FilterAuids rarity check had the greatest reduction to pipeline outputs (121 K ASVs removed and 6.0 % of  
 448 reads), followed by the length filter (4 K ASVs and 2.7 % of reads; Tables 3 and 5).  
 449



450

451

452 **Figure 4: Study-specific loss of reads at each stage of the post-pipeline workflow.** For each study the violin plots show how many  
 453 reads from the pipeline were removed by GatherAsvs due to length, the four filtering steps of FilterAuids, or WorkspaceStartup due to the

454 ASV having no annotation (shown in Fig. 1). Losses for all samples combined are shown in the box plot (top right). Bracketed numbers  
 455 after each study ID indicate the percentage of reads contributed to the *nifH* ASV database, e.g. 23.7 % of all the reads in the database were  
 456 from Hallstrom\_2021.

457

458

459 **Table 5. Quality filtering by the post-pipeline workflow.** For each study are shown the mean numbers of reads per sample that were  
 460 output by the DADA2 *nifH* pipeline and retained by the GatherAsvs, FilterAuids, and WorkspaceStartup stages of the post-pipeline  
 461 workflow. The Retained (%) column has the mean percentages of reads retained per sample (relative to column DADA2 pipeline values).  
 462 Additionally, the last three rows show the overall means, medians, and sums of reads across all samples and studies. Superscripts  
 463 correspond to stage numbers in Fig. 1 Post-pipeline stages. The GatherAsvs<sup>1</sup> column mainly reflects length filtering (200–450 nt), and the  
 464 WorkspaceStartup<sup>6</sup> column reflects discarding of ASVs that had no annotation. At each stage (column) the calculations include only the  
 465 samples that have >0 reads.

466

Study ID	DADA2 pipeline	Gather Asvs <sup>1</sup>	FilterAuids <sup>2</sup>				Workspace Startup <sup>6</sup>	Retained (%)	
			Small	Rare	Non-NifH	Length			
AK2HI	3.3E+04	3.3E+04	3.3E+04	3.3E+04	3.2E+04	3.0E+04	2.9E+04	89.2	
BentzonTilia_2015	4.1E+03	4.1E+03	4.0E+03	3.1E+03	3.1E+03	3.1E+03	3.0E+03	72.8	
Ding_2021	4.9E+04	4.9E+04	4.9E+04	4.6E+04	4.6E+04	4.5E+04	4.5E+04	92.2	
Gradoville_2020	2.4E+04	2.3E+04	2.3E+04	2.2E+04	2.1E+04	2.1E+04	2.0E+04	82.6	
Hallstrom_2021	1.4E+05	1.4E+05	1.4E+05	1.3E+05	1.3E+05	1.2E+05	1.2E+05	92.2	
Hallstrom_2022	6.6E+04	6.5E+04	6.5E+04	6.4E+04	6.4E+04	6.2E+04	6.2E+04	68.1	
Harding_2018	2.6E+04	2.6E+04	2.6E+04	2.4E+04	2.3E+04	2.0E+04	1.7E+04	75.6	
Mulholland_2018	1.0E+05	1.0E+05	1.0E+05	9.5E+04	9.3E+04	8.8E+04	8.4E+04	80.0	
NEMO	3.3E+04	3.3E+04	3.3E+04	3.2E+04	3.2E+04	3.0E+04	3.0E+04	84.2	
Raes_2020	6.5E+04	6.5E+04	6.5E+04	6.1E+04	6.1E+04	6.0E+04	5.9E+04	75.3	
Sato_2021	2.9E+04	2.7E+04	2.7E+04	2.2E+04	2.0E+04	1.5E+04	1.4E+04	49.2	
Selden_2021	8.0E+04	8.0E+04	8.0E+04	6.0E+04	5.2E+04	4.9E+04	4.5E+04	59.0	
Shiozaki_2017	1.6E+04	1.6E+04	1.6E+04	1.5E+04	1.5E+04	1.4E+04	1.4E+04	82.5	
Shiozaki_2018GBC	2.2E+04	2.2E+04	2.2E+04	2.1E+04	2.1E+04	2.1E+04	2.1E+04	90.4	
Shiozaki_2018LNO	4.8E+04	4.8E+04	4.8E+04	4.6E+04	4.6E+04	4.6E+04	4.6E+04	95.0	
Shiozaki_2020	1.4E+05	1.4E+05	1.4E+05	1.4E+05	1.4E+05	1.4E+05	1.4E+05	76.6	
Tang_2020	3.4E+04	3.4E+04	3.4E+04	3.3E+04	3.3E+04	3.3E+04	3.3E+04	97.9	
TurkKubo_2021	3.6E+04	3.5E+04	3.5E+04	3.5E+04	3.5E+04	3.4E+04	3.3E+04	94.1	
Wu_2019	4.2E+04	4.2E+04	4.2E+04	4.1E+04	4.1E+04	4.1E+04	4.1E+04	96.3	
Wu_2021	5.2E+04	5.2E+04	5.2E+04	4.8E+04	4.8E+04	4.8E+04	4.8E+04	93.2	
All samples	mean	5.0E+04	4.9E+04	4.9E+04	4.6E+04	4.6E+04	4.5E+04	4.4E+04	80.9
	median	3.0E+04	3.0E+04	3.0E+04	2.9E+04	2.8E+04	2.7E+04	2.6E+04	93.0

and studies	sum	4.9E+07	4.8E+07	4.8E+07	4.6E+07	4.5E+07	4.4E+07	4.3E+07	90.0
-------------	-----	---------	---------	---------	---------	---------	---------	---------	------

467

468

469 Finally, ASVs were removed if they were classified as non-*nifH*, based on a strong alignment to sequences in NCBI nr that  
470 ARBitrator (Heller et al., 2014) classified as non-*nifH*. Specifically, an ASV was classified as non-*nifH* if the ratio of  
471 E-values for its best positive and negative hits, among sequences classified by ARBitrator, was >10. A total of 137,366 of the  
472 139,355 non-chimera ASVs had database hits which resulted in 50,233 positive, 20,528 negative, and 66,605 uncertain  
473 classifications. This approach was used to leverage ARBitrator's high specificity for detecting *nifH* as well as to enable users  
474 to identify ASVs that have high percent identity matches to sequences in GenBank. An alternative approach would have  
475 been to classify the ASVs based on their alignments to HMMs for NifH versus NifH-like proteins (e.g. protochlorophyllide  
476 reductase), used by the NifMAP pipeline for *nifH* operational taxonomic units (Angel et al., 2018). Finally, FilterAuids  
477 removed ASVs with lengths outside 281–359 nt, a total of 4338 ASVs comprising 1.2 million reads (Figs. 1, 4 and Tables 3  
478 and 5). After FilterAUIDs, the total number of samples in the dataset was reduced from 982 to 951 and the number of ASVs  
479 from 139,355 to 11,915.

480

481 FilterAuids also flagged a total of 2342 ASVs as possible PCR contaminants. Although we opted to flag, not remove, these  
482 ASVs, the workflow can be easily altered to remove contaminants. Most studies contained low levels of contamination ( $\leq 1$   
483 %) based on our criteria. However, several studies were flagged with  $\sim 9$ –29 % of their reads being similar to known  
484 contaminants. Identifying potential contaminants is challenging given their numerous sources, study specific nature (Zehr et  
485 al., 2003), and lack of control sequence data from blanks.

486

487 Next, AnnotateAuids assigned annotations using our three *nifH* reference databases and CART (Fig. 1). In total 9406 of the  
488 11,915 quality filtered ASVs were annotated, usually with multiple references (Fig. A1). Most (9322 ASVs) had hits to both  
489 genome879 and ARB2017, likely because the 879 sequenced diazotrophs had *nifH* homologs in GenBank that were found by  
490 ARBitrator. Fewer ASVs had hits to the databases that targeted UCYN-A oligos (217 ASVs) and other marine diazotrophs  
491 (938 ASVs; 211 ASVs also had UCYN-A hits). Most ASVs (9380 total) were assigned to NifH clusters 1–4 by CART  
492 (respectively, 4923; 101; 4205; and 151 ASVs), including five ASVs that had no hits to our databases. The majority of ASVs  
493 (9257 total) had open reading frames (ORFs) that contained paired cysteines and AMP which might coordinate the 4Fe-4S  
494 cluster, and all 9257 also had annotation from the reference databases or CART. A few ASVs had annotations but lacked  
495 residues to coordinate 4Fe-4S: 29 ORFs lacked the paired cysteines and another 120 ORFs had paired cysteines but not  
496 AMP, usually due to a substitution for M. The last step of AnnotateAuids assigned primary IDs (described above) to 9383  
497 ASVs. All of them were retained in the final stage of the post-pipeline workflow, WorkspaceStartup (below).

498

499 In the CMAP stage, sample collection metadata (date, latitude, longitude, and depth) were used to download CMAP  
500 environmental data (100 variables) for each sample in the *nifH* ASV database (Fig. 1). The CMAP data will enable analyses  
501 of potential factors that influence the global distribution of the diazotrophic community. Aggregated metadata for all samples  
502 are available in the *nifH* ASV database (metaTab.csv for sample metadata and cmapTab.csv for environmental data).

503

504 The last stage of the post-pipeline workflow is WorkspaceStartup, which generates the *nifH* ASV database (Fig. 1). ASVs  
505 with no annotation are removed as well as samples with zero total reads due to ASV filtering steps. The *nifH* ASV database  
506 consisted of 21 studies, 944 samples, 9383 AVS and 43.0 million total reads (Tables 3 and 5). The database is heavily biased  
507 toward euphotic zone DNA samples, with euphotic heuristically defined as follows: Samples were classified as coastal (<  
508 200 km from a major landmass) or open ocean. Euphotic samples were then identified as those collected above a depth cut  
509 off, 50 m for coastal samples and 100 m for open ocean. Samples obtained from DNA (n=836) far exceeded those from RNA  
510 (n=108) extracts. Likewise, a majority of the samples were from the euphotic zone (861 compared to 83 from the aphotic  
511 zone). The database also includes replicate samples (n=286) and size fractionated samples (n=170).

## 512 3.2 Global *nifH* ASV database

### 513 3.2.1. Comparison to an OTU database

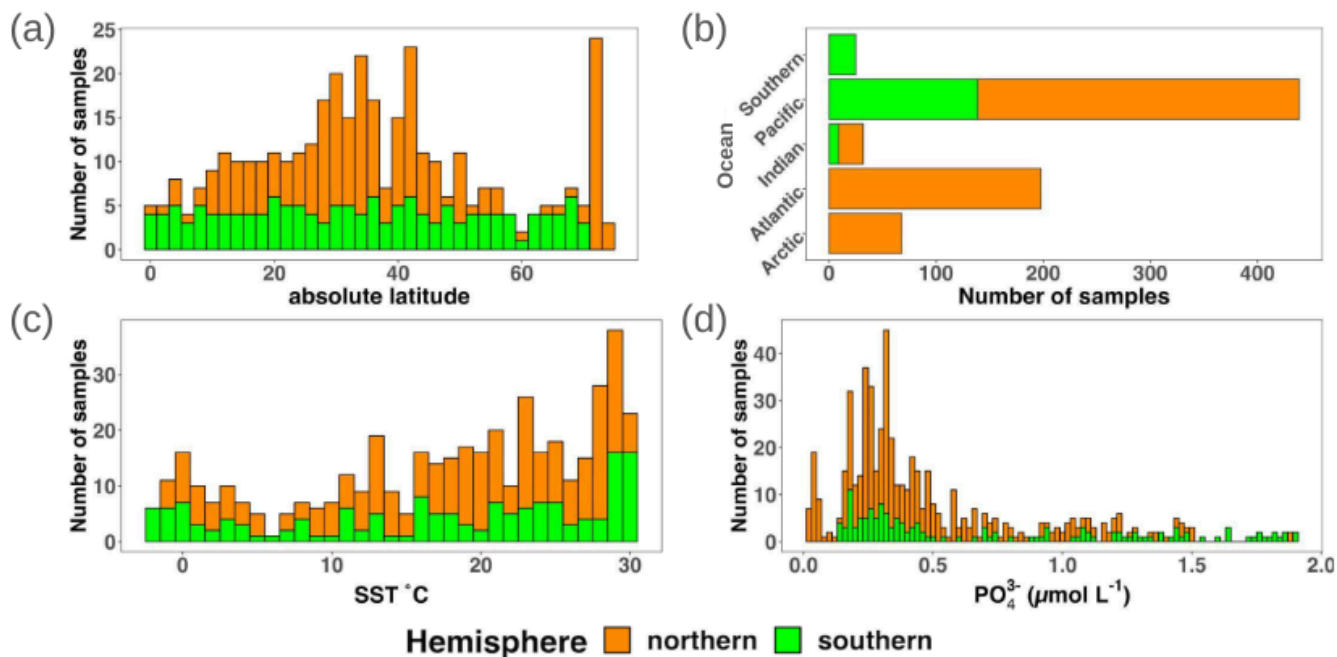
514 New studies with Illumina amplicon data have mainly used DADA2 (Callahan et al., 2016) and other methods that  
515 distinguish fine-scale variation from sequencing errors (Eren et al., 2014; Edgar, 2016b; Amir et al., 2017). Earlier studies,  
516 including 13 of the 19 previously published studies in the *nifH* ASV database (Table C1), used *de novo* operational  
517 taxonomic units (OTUs) which were obtained by clustering the sequences at 97 % nucleotide identity. OTUs masked  
518 sequencing errors as well as fine-scale variation and had other disadvantages compared to ASV approaches (Callahan et al.,  
519 2017). Although cross-study comparisons ideally will use the same pipeline for all the studies—the motivation for our  
520 workflow—previously published results should be considered. Therefore, for each study in the *nifH* ASV database,  
521 diazotroph communities were compared to versions generated using the NifMAP OTU pipeline (Appendix C). The ASV and  
522 OTU communities mainly had similar *nifH* clusters, except for several studies where the workflow retained substantially  
523 more sequencing reads (Fig. C1, Table C1).

### 524 3.2.2. Sample Distribution

525 Investigations of N<sub>2</sub> fixation and diazotrophic communities have focused on specific ocean regions and this is reflected by  
526 the uneven global distribution of *nifH* amplicon datasets in the *nifH* ASV database (Figs. 2, 5a, and 5b). There is an outsized  
527 influence of the northern hemisphere, especially in the Pacific Ocean where most of the database samples were located (439)  
528 and 68.3 % of these samples originated from the northern hemisphere (Figs. 2, 5a, 5b, and 6). Ten studies are found within  
529 the Pacific, with several containing >50 samples (Figs. 2 and 6). Notably, Raes\_2020 (Raes et al., 2020) is the largest dataset

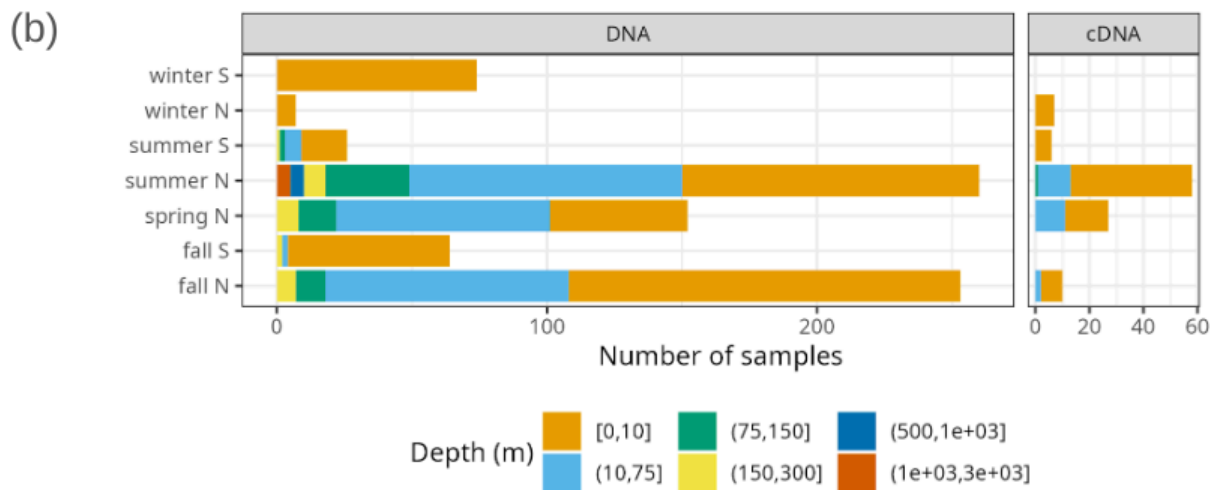
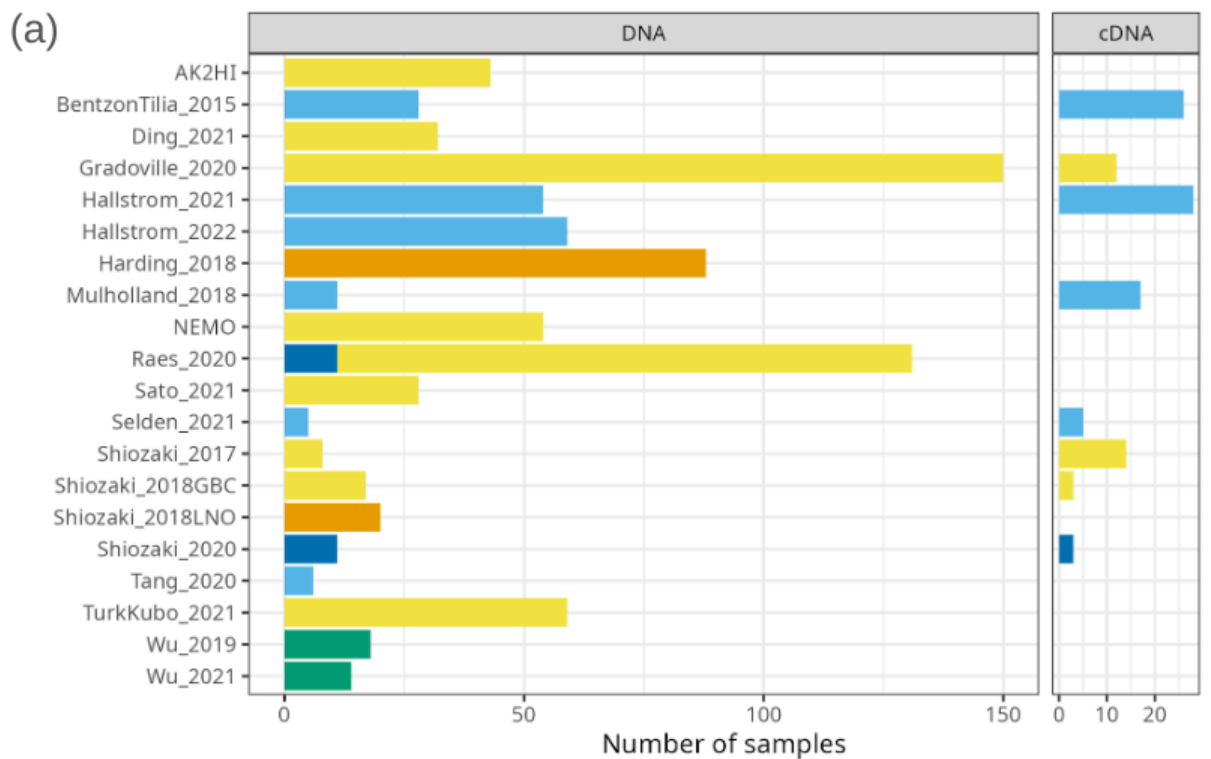


530 stretching from the equator to the Southern Ocean, making up almost the entirety of the southern hemisphere Pacific samples  
 531 (Figs. 2 and 6). Two new studies carried out in the North Pacific constitute the only previously unpublished data of the *nifH*  
 532 ASV database (Table 1). AK2HI was a latitudinal transect from Alaska (U.S.) to Hawaii (U.S.) and NEMO was a  
 533 longitudinal transect across the Eastern North Pacific from San Diego, CA (U.S.) to Hawaii (U.S.) (Fig. 2; Sect. 2.2.2). The  
 534 amplicon data compiled for the *nifH* ASV database was primarily generated from DNA, with most RNA samples deriving  
 535 from Atlantic Ocean studies and no contribution from RNA samples in the Arctic or Indian Oceans (Fig. 6).  
 536



537  
 538 **Figure 5. Location, temperature, and phosphate distributions of the *nifH* ASV database.** The number of samples from the *nifH* ASV  
 539 database by (a) absolute latitude, (b) the world's oceans, (c) sea surface temperature (SST, °C) and (d) Pisces-derived PO<sub>4</sub><sup>3-</sup> (µmol L<sup>-1</sup>).  
 540 Environmental data, (c) and (d), were retrieved from the CMAP data portal. All bars are stacked.

541



542

543

544 **Figure 6. Samples in the *nifH* ASV database by collection location, season, and amplicon type.** The number of samples from each  
 545 study are shown by ocean and study (a), and by the collection season, hemisphere, and depth (b). For both panels the amplicon type (DNA  
 546 or cDNA) is shown, but x axis scales differ between (a) and (b). See Table 1 for citations for the studies in (a).

547

548

549 Under-sampled regions include the Eastern South Pacific (n=6) and the Western Indian Ocean (n=0) (Figs. 2, 5a, and 6a).  
550 Only two studies originated from the Indian Ocean, a unique environment with intense weather and shifting circulation  
551 patterns that include monsoon seasons and upwelling conditions that will require much greater sampling coverage to capture  
552 diazotroph biogeography. No South Atlantic samples were found during compilation that met the criteria for inclusion in the  
553 *nifH* ASV database, though there are several studies from this region (Table A1). Most Atlantic Ocean samples were coastal  
554 and from the North Atlantic. Thus, the Atlantic subtropical gyres, which are known to host diverse diazotrophs (Langlois et  
555 al., 2005), are underrepresented by *nifH* amplicon data (Fig. 2).

556

557 Tropical and subtropical regions, often associated with high temperatures and low nutrients, are highly represented in the  
558 database (Figs. 2 and 5a). This likely influenced the ranges of environmental variables with most samples in the database  
559 originating from locations with SST above 15 °C and PO<sub>4</sub><sup>3-</sup> below 0.5 μmol L<sup>-1</sup> (Figs. 5c and 5d). Northern hemisphere  
560 samples were collected in all seasons, though fewer from the winter. In contrast, most southern hemisphere samples were  
561 collected in the winter and fall (Fig. 6b). While most DNA samples are from the euphotic zone (Fig. 6b), cDNA samples are  
562 almost exclusively from the euphotic zone, and mainly from the northern hemisphere during the spring and summer (Fig.  
563 6b), indicating an incomplete picture of diazotroph activity.

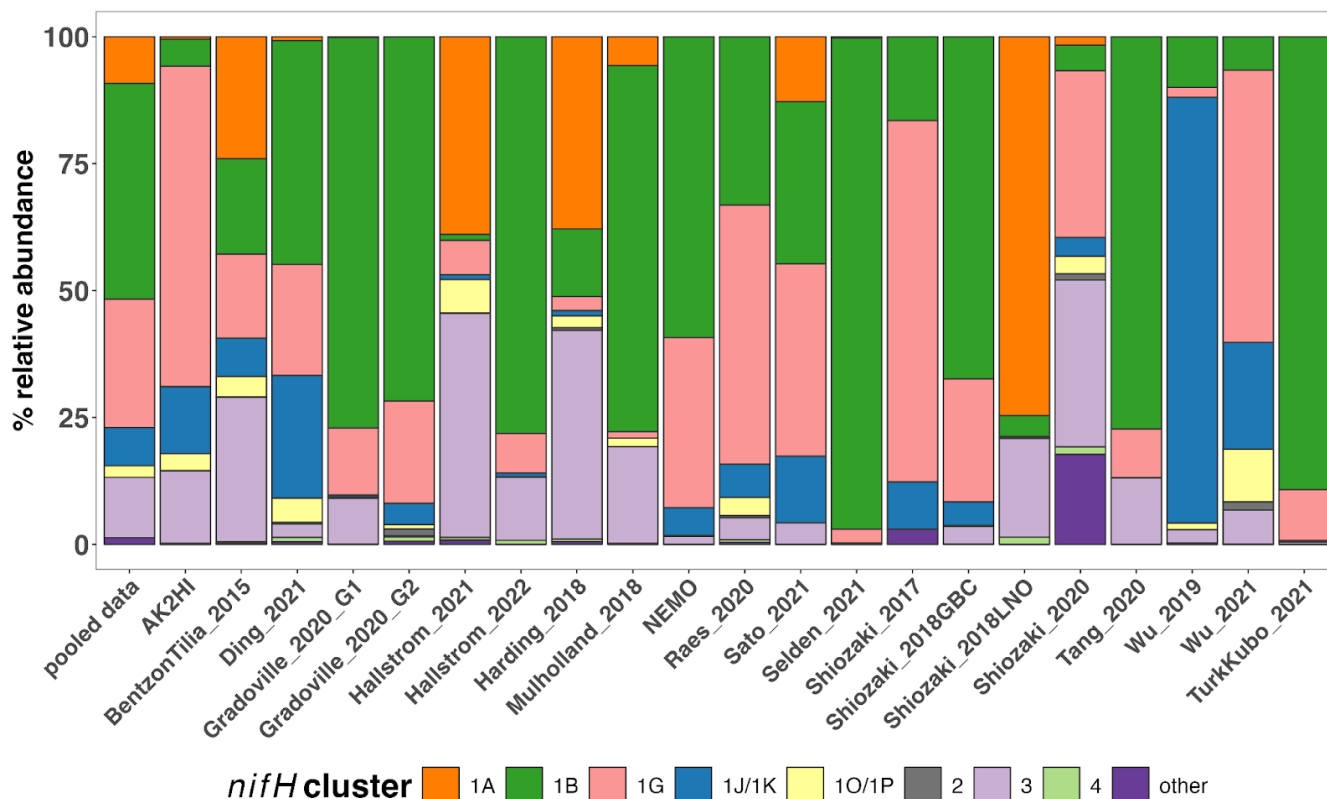
564

565 The disproportionate spatial and seasonal coverage between hemispheres in the *nifH* ASV database mirrors collection biases  
566 in other N<sub>2</sub> fixation metrics including: N<sub>2</sub> fixation rate measurements; diazotroph cell counts; and *nifH* qPCR data, which are  
567 heavily sourced from the North Atlantic (Shao et al., 2023) or, when targeting NCDs, also the North Pacific (Turk-Kubo et  
568 al., 2022). These biases underscore the need for future work in understudied regions and seasons.

### 569 3.3 Study-specific patterns in global diazotroph assemblages in the DNA dataset

570 To demonstrate how the *nifH* ASV database can be used, a subset of the data was created that comprised all DNA samples  
571 (88.8 % of the total dataset; Fig. 7) and referred to herein as the “DNA dataset.” Samples derived from cDNA (n=108; Fig.  
572 6) were removed. Replicate samples (n=286) or those with multiple size fractions (n=170) were combined by averaging  
573 across replicates or size fractions. This reduced the number of DNA samples to 762 and the total number of reads in the  
574 count table to 36.6 million from 43.0 million.

575



576

577 **Figure 7. Study-specific diazotroph assemblage patterns in the DNA dataset.** The percentage of relative abundance over the DNA  
 578 dataset for each major *nifH* cluster. The first column ('pooled data') uses all the compiled data while each subsequent column only uses  
 579 data from the indicated study. Colors represent different *nifH* subclusters; 'other' are the remaining *nifH* clusters.

580

581

582 As demonstrated in a previous global analysis of diazotroph assemblages (Farnelid et al., 2011), cyanobacterial sequences  
 583 (cluster 1B) dominate the samples, making up 42 % of the total relative abundance (Fig. 7). Though photosynthetic  
 584 cyanobacteria would be expected to thrive in euphotic waters, NCDs are also widespread in the ocean surface (Langlois et  
 585 al., 2005; Delmont et al., 2018; Delmont et al., 2022; Pierella Karlusich et al., 2021; Turk-Kubo et al., 2022). Indeed, among  
 586 the NCDs,  $\gamma$ -proteobacteria (*nifH* cluster 1G) were the most prevalent, comprising 27 % of the total relative abundance,  
 587 while  $\delta$ -proteobacteria (clusters 1A and 3) accounted for 21 % of the total relative abundance of the DNA dataset (Fig. 7).  
 588 Less prominent clusters 1J/1K ( $\alpha$ - and  $\beta$ -proteobacteria) and 1O/1P ( $\gamma$ -/ $\beta$ -proteobacteria and Deferribacteres) were 4 % and 3  
 589 % of the relative abundance, respectively. The remaining ASVs comprised <1.5 % of the total relative abundance and came  
 590 from clusters associated with nitrogenases that do not use iron (e.g. cluster 2) or that are uncharacterized (cluster 4) (Fig. 7).

591

592 Cluster 1B (cyanobacteria) were generally high in individual studies across the *nifH* DNA dataset, comprising  $\geq 25$  % of the  
 593 community in two-thirds of the studies (Fig. 7), which is the highest of any cluster. Studies carried out in polar regions

594 (Harding\_2018, Shiozaki\_2018LNO, Shiozaki\_2020) and the Indian Ocean (Wu\_2019 and Wu\_2021) were distinct from  
595 this pattern, with low relative abundances of cluster 1B. Instead, Arctic studies had high relative abundances of cluster 1A  
596 and 3 (both primarily comprised of  $\delta$ -proteobacteria) and while clusters 1J/1K ( $\alpha$ - and  $\beta$ -proteobacteria) and 1O/1P  
597 ( $\gamma$ - $\beta$ -proteobacteria and Deferribacteres) were the predominant groups in the Indian Ocean.

598

599 The second most abundant group was the cluster 1G ( $\gamma$ -proteobacteria), making up ca. 25 % of the total relative abundance  
600 across the DNA dataset, with study-specific relative abundances greater than 25 % in eight out of 21 studies (Fig. 7).  
601 Members of this group were often found at high relative abundances in Pacific Ocean studies (AK2HI, NEMO, Raes\_2020,  
602 Sato\_2021, Shiozaki\_2017), as well as in other ocean regions including the Atlantic (BentzonTillia\_2015), Indian (Wu\_2021)  
603 and Southern Ocean (Shiozaki\_2020). The notable exception is in Arctic studies (Harding\_2018, Shiozaki\_2018LNO) where  
604 cluster 1G was almost absent (Fig. 7).

605

606 In several studies, including BentzonTillia\_2015, Hallstrom\_2021, Mulholland\_2018, Selden\_2021, Tang\_2020, and  
607 Hallstrom\_2022, diazotroph assemblages had high relative abundances of putative  $\delta$ -proteobacteria (clusters 1A and 3),  
608 reflecting possibly a coastal/shelf or upwelling signature (Figs. 2 and 7). The only study with samples primarily from the  
609 Southern Ocean (Shiozaki\_2020) was also the only study with a large portion of *nifH* cluster 1E (*Bacillota*).

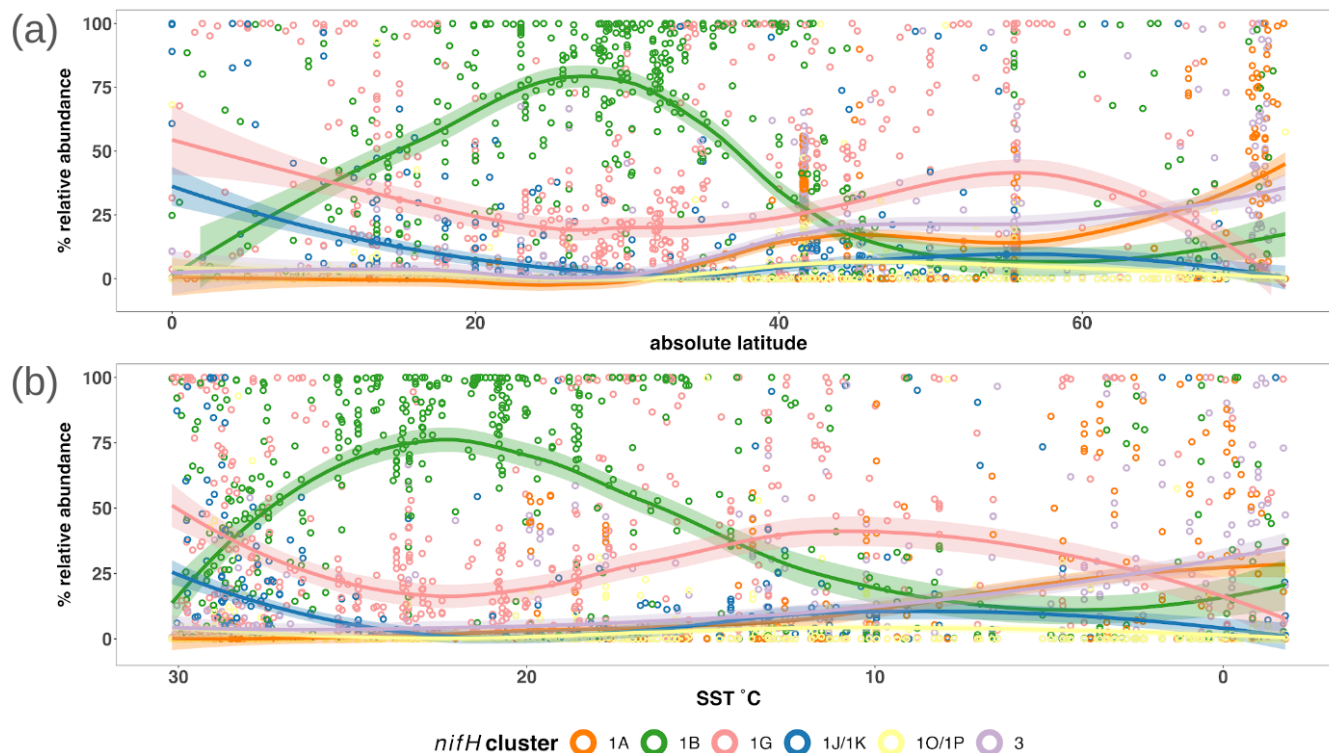
### 610 3.3.2 Emerging patterns in global diazotroph assemblages across the DNA dataset

611 The *nifH* ASV database enables new analyses of global diazotroph biogeography in the context of environmental parameters,  
612 through co-localization with satellite and model outputs publicly available through CMAP (Ashkezari et al., 2021). To  
613 demonstrate the utility of the *nifH* ASV database, we present here patterns in relative abundances of *nifH* clusters across  
614 absolute latitude and SST in the DNA dataset. Cosmopolitan distributions were evident for  $\gamma$ -proteobacterial (1G) and  
615 cyanobacterial diazotrophs (1B; Fig. 8a), corroborating and extending previous findings (Farnelid et al., 2011; Shao and Luo,  
616 2022; Halm et al., 2012; Fernandez et al., 2011; Löscher et al., 2014; Cheung et al., 2016). At low to mid latitudes,  
617  $\gamma$ -proteobacterial (1G) diazotrophs generally had high relative abundances and were often the dominant taxa when present.  
618 However, they declined within the gyre regions, ranging between ~25–50 % of the population when present, while  
619 cyanobacterial diazotrophs (1B) increased and became dominant in the subtropical gyres (Fig. 8a). Notably, cluster 1G  
620 diazotrophs reached high relative abundances in each transitional zone, before mainly disappearing at latitudes above 56°  
621 (Fig. 8a). However, as mentioned previously, sampling bias likely plays a large role at these higher latitudes where the  
622 number of studies and samples are sparse (Figs. 2 and 5).

623

624 Clusters 1B and 1G were both detected over the full range of SST (approximately -2–30 °C) but peaks in their relative  
625 abundances occurred in distinct SST ranges. Cyanobacterial diazotrophs had multiple peaks in relative abundance in waters  
626 >18 °C underscoring their dominance in tropical gyre regions (Fig. 8b). The 1G cluster also spanned the entire temperature

627 spectrum but had notably higher presence and relative abundance above SSTs of 8 °C and 11 °C, respectively (Fig. 8b). The  
628 overlap between 1G and 1B has been reported previously, however the factors controlling this are unknown (Moisander et  
629 al., 2014; Shiozaki et al., 2017; Shiozaki et al., 2018b; Liu et al., 2020; Tang et al., 2020; Messer et al., 2015).  
630



631

632 **Figure 8: Influence of SST on the global distribution of major *nifH* clusters in the photic zone of the DNA dataset.** The relative  
633 abundance of *nifH* genes for each major *nifH* cluster from every photic zone sample compiled in the DNA dataset versus (a) absolute  
634 latitude and (b) SST. Smoothing averages (lines) were calculated using local polynomial regression fitting (LOESS) with 95 % confidence  
635 intervals (translucent colored areas). Each color represents a different *nifH* cluster. SST in (b) is from warmest to coldest temperatures to  
636 show that trends are similar to those in (a).

637

638  $\delta$ -proteobacterial diazotrophs (clusters 1A and 3) were generally found in cooler, higher latitude waters. Notably, both  
639 clusters 1A and 3 were mainly found below  $\sim 10^{\circ}\text{C}$  (Fig. 8b).  $\delta$ -proteobacteria associated with cluster 1A were generally  
640 found at latitudes  $>32^{\circ}$  and reached maximum relative abundances near the poles, including in the Beaufort Sea, the highest  
641 latitude region surveyed ( $72^{\circ}$ ; Figs. 2, 5, and 8a). The vast majority of cluster 1A  $\delta$ -proteobacteria were found at SST  $\leq 5^{\circ}\text{C}$   
642 (Fig. 8b). Though cluster 3 and 1A distributions were similar, cluster 3 showed broader spatial and temperature ranges, with  
643 consistent but low relative abundances in the subtropics and tropics (Fig. 8).

644

645 In contrast, the relative abundances of cluster 1J/1K and 1O/1P diazotrophs declined as SST decreased and latitude  
646 increased, becoming rare at higher latitudes (Fig 8). The highest relative abundances for these clusters were observed near  
647 the equator, and in some cases, comprised 100% of the diazotroph assemblage in high SST, tropical samples. These patterns  
648 suggest that temperature was an important factor controlling the narrow SST band ( $\geq 26$  °C) clusters 1J/1K and 1O/1P  
649 occupied, establishing them as the *nifH* clusters with the smallest geographic range in the *nifH* ASV database (Fig. 8).

650

651

### 652 3.4 Limits and caveats to interpreting *nifH* amplicon data

653 The PCR amplification of the *nifH* gene and its transcripts has been vital in advancing the knowledge of diazotroph ecology  
654 due to its high sensitivity, detecting diazotrophs at abundances that are often orders of magnitude lower than other marine  
655 microbes. This approach has facilitated the discovery of many novel diazotrophs, and provided the first evidence of the  
656 widespread distribution of unicellular diazotrophs throughout the open oceans (Falcón et al., 2004; Falcón et al., 2002; Zehr  
657 et al., 1998; Zehr et al., 2001). Advances in HTS technologies have revealed diverse diazotrophic assemblages, including the  
658 ubiquitously distributed NCDs (Turk-Kubo et al., 2014; Shiozaki et al., 2017; Raes et al., 2020). These discoveries have  
659 fostered a new perspective of global diazotrophic ecology (Zehr and Capone, 2020), improved our models of diazotrophic  
660 distributions and global N fixation rates (Tang et al., 2019) and will continue to drive new research questions.

661

662 However, interpreting *nifH* PCR-based data requires the consideration of several important caveats. Diazotrophs constitute a  
663 small fraction of the total microbial community, and thus often require numerous PCR cycles in conjunction with nested  
664 PCR for detection. Increasing the number of cycles can exacerbate known amplification biases (Turk et al., 2011) and  
665 increase the likelihood of detecting contaminant sequences (Zehr et al., 2003). Strategies to mitigate and assess  
666 contamination exist, e.g., by employing ultrafiltration of reagents and including blanks at different stages of the sampling and  
667 sequencing process (Bostrom et al., 2007; Farnelid et al., 2011; Blais et al., 2012; Moisander et al., 2014; Langlois et al.,  
668 2015; Fernandez-Mendez et al., 2016; Cheung et al., 2021), but such strategies have not been universally adopted.  
669 Additionally, relative abundances of PCR amplicons cannot easily be related to absolute abundances. For example, the  
670 relative abundance of a taxon can change even if its absolute abundance remains constant, or the relative abundance can  
671 remain constant despite changes in the total assemblage size. Moreover, the complexity of the diazotroph assemblage can, if  
672 the HTS sequencing depth is insufficient, cause rare ASVs to go undetected, or have relative abundances which are too low  
673 to interpret.

674

675 Primary objectives in studying marine diazotrophic populations include understanding the contribution of each group to N<sub>2</sub>  
676 fixation, the factors influencing their activity, and their global distributions. The relative abundances of *nifH* genes and  
677 transcripts estimated by the workflow can point to potentially significant contributors to N<sub>2</sub> fixation rates. Yet, the presence



678 of *nifH* genes or transcripts does not always correlate with N<sub>2</sub> fixation rates (e.g. Gradoville et al., 2017). This underscores  
679 the need for cell-specific rates to better constrain N<sub>2</sub> fixation, the assemblages driving given rates, and the taxa-specific  
680 regulatory factors of N<sub>2</sub> fixation to better constrain global biogeochemical modeling.

681

682 Various methods are available to target specific diazotroph taxa over space and time (e.g. qPCR/ddPCR, fluorescent in situ  
683 hybridization (FISH)-based methods). Universal PCR assays, e.g., those used in the studies compiled here (*nifH1-4*), are an  
684 important complement because they better capture the overall diversity of the diazotrophic assemblage. Unlike primers  
685 designed for specific sequences, universal primers can amplify unknown or ambiguous sequences, enabling the discovery of  
686 genetic diversity. This includes microdiversity, where sequences show subtle variations from known ones, or even  
687 identifying entirely novel taxa. Primers specific to novel sequences can then be developed for use in the mentioned  
688 quantitative methods, enabling experiments to characterize the growth, activity, and controlling factors/dynamics of putative  
689 diazotrophs growth.

690

691 Tools like RT-qPCR, where transcript abundances are assessed directly, or FISH-based methods where single-cells are  
692 identified for cell-specific analysis, provide complementary perspectives into the activities of putative diazotrophs.  
693 Enumerating diazotrophs using techniques like these can help standardize the relative abundances associated with amplicon  
694 sequencing via matching taxa across each method. By assessing diversity and abundance simultaneously, major players can  
695 potentially be identified and monitored.

696

697 Through genome reconstruction, `omics studies can enhance the characterization of putative diazotroph amplicon sequences  
698 by providing a robust suite of associated genetic data, e.g., taxonomic, phylogenetic, and metabolic. Previous studies have  
699 led to the assembly of dozens of diazotrophic genomes (Delmont et al., 2022; Delmont et al., 2018). However, `omics  
700 methods often require massive amounts of data to detect rare community members, and linking genes of interest to other  
701 genomic information, e.g., taxonomy, remains quite difficult. Gene-specific models are also required to retrieve diazotrophic  
702 information and these models can benefit greatly from the high quality diazotrophic sequences of the *nifH* ASV database. In  
703 summary, the complementary perspectives afforded by the methods just described should all be used to obtain robust insights  
704 into diazotrophic assemblages.

705

#### 706 **4 Data availability**

707 The *nifH* ASV database is freely available in Figshare (<https://doi.org/10.6084/m9.figshare.23795943.v2>; Morando et al.,  
708 2024a). HTS datasets for the 21 studies in the database can be obtained from the NCBI Sequence Read Archive using the  
709 NCBI BioProject accessions in Table 1.

## 710 5 Code availability

711 The workflow used to generate the *nifH* ASV database is freely available in two GitHub repositories, one for the DADA2  
712 *nifH* pipeline ([https://github.com/jdmagasin/nifH\\_amplicons\\_DADA2](https://github.com/jdmagasin/nifH_amplicons_DADA2); Morando et al., 2024b) and one for the post-pipeline  
713 stages (<https://github.com/jdmagasin/nifH-ASV-workflow>; Morando et al., 2024c).

## 714 6 Conclusions

715 The workflow and *nifH* ASV database represent a significant step towards a unified framework that facilitates cross-study  
716 comparisons of marine diazotroph diversity and biogeography. Furthermore, they could guide future research, including  
717 cruise planning, e.g., focusing more on the southern hemisphere and areas outside of the tropics, and molecular assay  
718 development, e.g., assays to characterize NCDs for single-cell activity rates.

719

720 To demonstrate the utility of our framework, the DNA dataset was used to identify potentially important ASVs and  
721 diazotrophic groups, establishing global biogeographic patterns from this aggregated amplicon data. Cyanobacteria were the  
722 dominant diazotrophic group, but cumulatively the NCDs made up more than half of the total data. Distinct latitudinal  
723 patterns were seen among these major diazotrophic groups, with NCDs (clusters 1G, 1J/K, 1O/1P, 1A, and 3) having a  
724 greater contribution to relative abundances near the equator and at higher latitudes, while cyanobacteria (1B) comprised a  
725 majority of the diazotroph assemblage in the subtropics. SST appeared to restrict and differentiate the biogeography of  
726 clusters 1J/1K and 1O/1P (warm tropics/subtropics) from clusters 3 and 1A (cool, high latitude waters), but did not play as  
727 large of a role for the biogeography of clusters 1B and 1G.

728

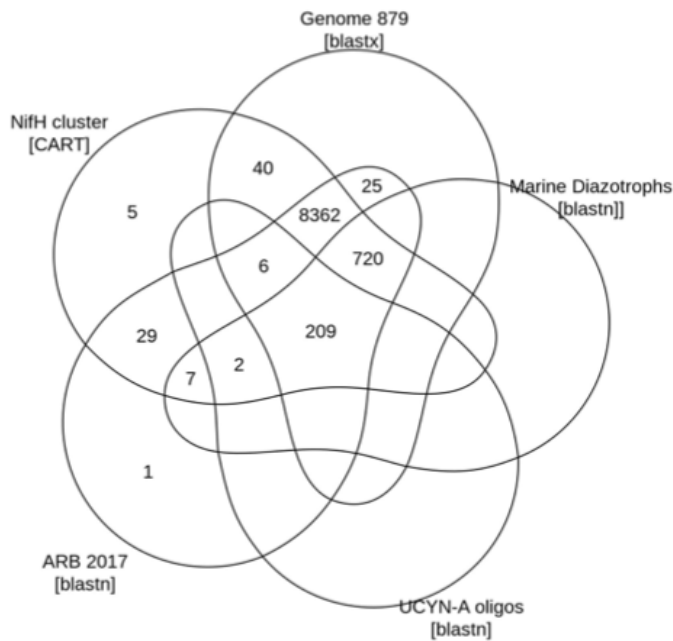
729 We provide the workflow and database for future investigations into the ecological factors driving global diazotrophic  
730 biogeography and responses to a changing climate. Ultimately, we hope that insights derived from the use of our framework  
731 will inform global biogeochemical models and improve predictions of future assemblages.

732

## 733 Appendix A:

734 Figures:

735



736

737 **Figure A1. ASV annotations.** The Venn diagram summarizes annotations assigned to 9406 ASVs during the AnnotateAuids stage of the  
 738 workflow (Fig. 1). Numbers indicate how many ASVs received each type of annotation. Of the 11,915 ASVs from the preceding  
 739 workflow stage, FilterAuids, only the 9406 ASVs shown received annotations.

740

741 Tables:

742 **Table A1. Compiled *nifH* amplicon studies.** Information on all studies compiled to generate the *nifH* ASV database, as well as studies  
 743 that were not ultimately included and the reasons for this. The table provides the study ID used to refer to each dataset, the NCBI  
 744 BioProject accession, the number of samples, and the DOI of the publication in which the dataset became public. \*: Data were obtained  
 745 from authors, not the SRA.

Study ID	Sam ples	NCBI BioProject	Reference	DOI	In <i>nifH</i> ASV database?
AK2HI	43	PRJNA1062410	This study	n/a	Yes
BentzonTilia_2015	56	PRJNA239310	Bentzon-Tilia et al., 2015	10.1038/ismej.2014.119	Yes
Cabello 2020	75	PRJNA605009	Cabello et al., 2020	10.1111/jpy.13045	No. Time series samples
Ding_2021	32	SUB7406573	Ding et al., 2021	10.3390/biology10060555	Yes
Farnelid 2019	155	PRJNA392595	Farnelid et al., 2019	10.1038/s41396-018-0259-x	No. Particle enrichment samples
Gérikas Ribeiro 2018	55	PRJNA377956	Gérikas Ribeiro et al., 2018	10.1038/s41396-018-0050-z	No. Samples had very few sequences
Gradoville 2017 Frontiers	45	PRJNA358796	Gradoville et al., 2017	10.3389/fmicb.2017.01122	No. Perturbation experiments
Gradoville_2020_G1	111	PRJNA530276	Gradoville et al., 2020	10.1002/ln0.11423	Yes
Gradoville_2020_G2	56	PRJNA530276	Gradoville et al., 2020	10.1002/ln0.11423	Yes

<b>Hallstrom_2021</b>	82	PRJNA656687	Hallstrøm et al., 2022b	10.1002/Ino.11997	Yes
<b>Hallstrom_2022</b>	83	PRJNA756869	Hallstrøm et al., 2022a	10.1007/s10533-022-00940-w	Yes
<b>Harding_2018</b>	91	PRJNA476143	Harding et al., 2018	10.1073/pnas.1813658115	Yes
<b>Li 2018</b>	16	PRJNA434503	Li et al., 2018	10.3389/fmicb.2018.00797	No. Issues merging reads
<b>Mulholland_2018</b>	29	PRJNA841982	Mulholland et al., 2019	10.1029/2018GB006130	Yes
<b>NEMO</b>	56	PRJNA1062391	This study	n/a	Yes
<b>Raes_2020</b>	121	PRJNA385736	Raes et al., 2020	10.3389/fmars.2020.00389	Yes
<b>Rahav 2016</b>	n/a	n/a	Rahav et al., 2016	10.1038/srep27858	No. Samples sorted prior to sequencing
<b>Sato_2021</b>	28	PRJDB10819	Sato et al., 2021	10.1029/2020JC017071	Yes
<b>Selden_2021</b>	10	PRJNA683637	Selden et al., 2021	10.1002/Ino.11727	Yes
<b>Shiozaki_2017*</b>	22	PRJDB5199	Shiozaki et al., 2017	10.1002/2017GB005681	Yes
<b>Shiozaki_2018GBC*</b>	20	PRJDB6603	Shiozaki et al., 2018b	10.1029/2017GB005869	Yes
<b>Shiozaki_2018LNO</b>	20	PRJDB5679	Shiozaki et al., 2018a	10.1002/Ino.10933	Yes
<b>Shiozaki_2020</b>	14	PRJDB9222	Shiozaki et al., 2020	10.1038/s41561-020-00651-7	Yes
<b>Tang_2020</b>	6	PRJNA554315	Tang et al., 2020	10.1038/s41396-020-0703-6	Yes
<b>Turk-Kubo 2015</b>	11	PRJNA300416	Turk-Kubo et al., 2015	10.5194/bg-12-7435-2015	No. Mesocosm samples
<b>TurkKubo_2021</b>	136	PRJNA695866	Turk-Kubo et al., 2021	10.1038/s43705-021-00039-7	Yes
<b>Wu_2019</b>	18	PRJNA438304	Wu et al., 2019	10.1007/s00248-019-01355-1	Yes
<b>Wu_2021*</b>	14	PRJNA637983	Wu et al., 2021	10.1007/s10021-021-00702-z	Yes

746

747

#### 748 **Appendix B: Read trimming method effects on workflow outputs**

749 It is well-established that error rates increase with the number of PCR cycles during Illumina sequencing (Manley et al.,  
750 2016). DADA2 trims the reads to remove the low-quality tails, an important early step that impacts the proportion of  
751 sequences retained during quality-filtering and merging, as well as the ASVs detected (Fig. 1). Usually sequencing quality  
752 plots are inspected to identify a trimming length that will on average cut the reads before quality declines significantly.  
753 However, inspecting tens to hundreds of quality plots (depending on the study size) is laborious and unsystematic. For the  
754 present work, the pipeline ancillary script estimateTrimLengths.R was used to efficiently identify lengths that maximized the  
755 percentages of reads retained for each study (Section 2.3.2). The optimized lengths appeared in the parameter files as  
756 truncLen.fwd and truncLen.rev used by DADA2 filterAndTrim (Table 2).

757

758 An alternative to fixed-length trimming is to trim each read based on its individual quality profile, at the first position where  
759 the estimated sequencing error rate exceeds a threshold specified in the truncQ parameter to filterAndTrim (Table 2). This

760 approach might reduce mismatches in the overlapping regions during the merge step and thus retain more read pairs.  
761 However, spurious low-quality bases could cause overly aggressive trimming, and picking a threshold that allows most  
762 sequences to overlap is not straightforward.

763

764 The quality of the raw sequencing data is a critical factor in the generation of the final ASV table. When analyzing a new  
765 dataset, testing both the fixed-length (`truncLen`) and quality-based (`truncQ`) trimming methods is suggested because they are  
766 fundamentally different and `filterAndTrim` impacts all downstream DADA2 steps. If both methods produce similar ASVs  
767 and abundances, additional parameter tuning is unlikely to impact the analysis meaningfully.

768

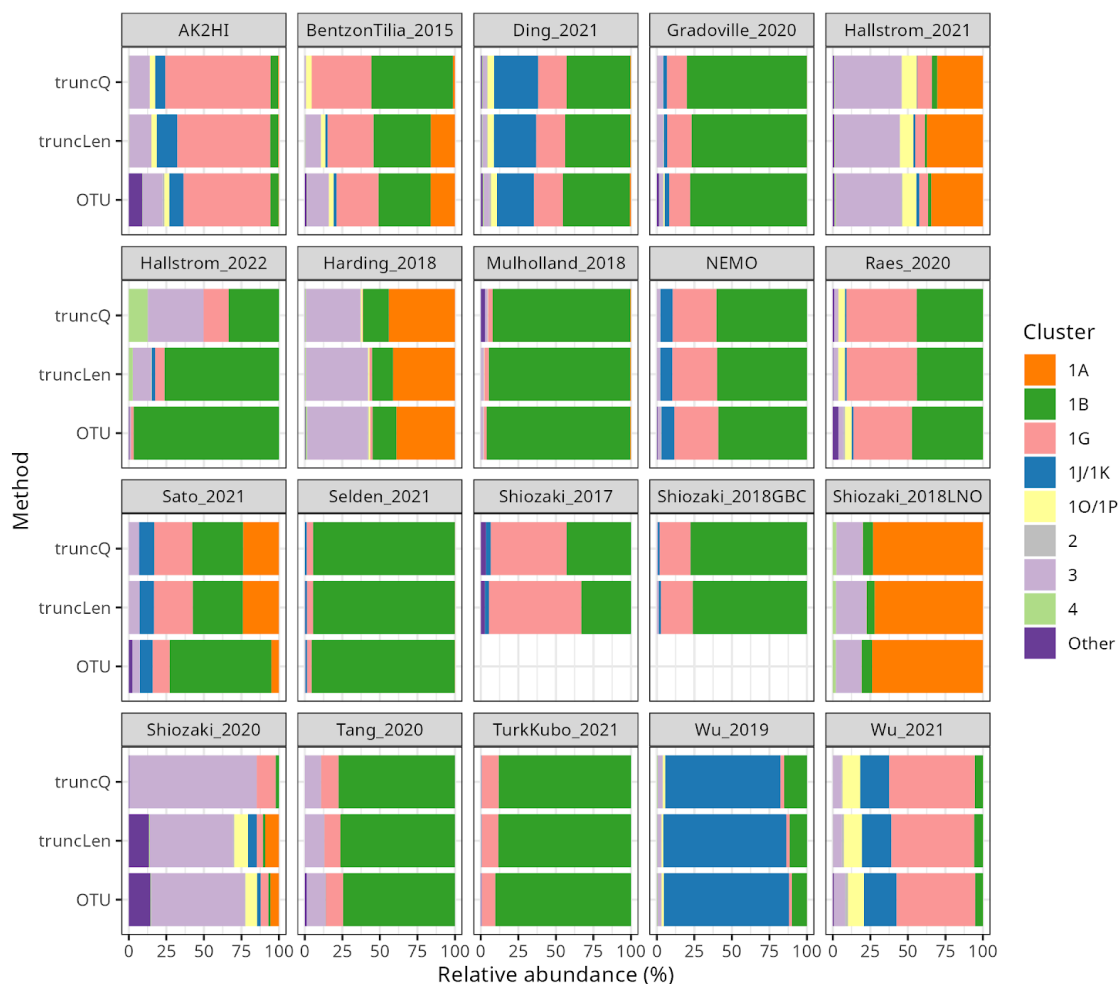
769 To illustrate how the trimming approach can impact workflow outputs, a version of the *nifH* ASV database was generated as  
770 shown in Figure 1 except that reads were trimmed at the first position where the estimated error rate was  $>2.5\%$  (`truncQ` =  
771 16 in Table 2). This threshold typically produces forward and reverse ASVs of sufficient length to overlap without  
772 mismatches. The `truncQ` version of the database had substantially fewer samples, reads, and ASVs (Table B1), partly  
773 because `truncQ` appeared more affected by low quality reads (discussed below). Only 1783 ASVs out of 9383 in the *nifH*  
774 ASV database were detected by both trimming methods, but they comprised 88.3 % of the total reads in the database (Table  
775 B1). The 7600 ASVs (16.7 % of reads) that were found only using `truncLen` had mainly low abundances and were detected  
776 mainly in one to several samples. Although `truncQ` was less sensitive to rare ASVs, for most studies the relative abundances  
777 of *nifH* groups were similar using either trimming approach (Fig. B1).

778

779 There were three exceptions where sequencing quality issues caused substantial differences in the results from `truncQ` and  
780 `truncLen`, BentzonTilia\_2015, Hallstrom\_2022, and Shiozaki\_2020. Using either trimming method, all three studies lost  
781 high percentages of reads during `filterAndTrim` (Fig. 3; losses using `truncQ` were comparable). This indicates that  
782 sequencing errors remained after trimming ( $>2$  errors in the trimmed forward reads and  $>4$  in the reverse; `maxEE` in Table 2).  
783 However, the subsequent losses during `mergePairs` were much higher using `truncQ` (vs. `truncLen`), respectively 58 % (10 %),  
784 61 % (5 %), and 72 % (6 %) of reads. This suggests that trimming with `truncQ=16` more frequently produced reads that  
785 failed to overlap during the merge step. For these three studies the workflow discarded many samples due to having  $\leq 500$   
786 reads, but more with `truncQ` (vs. `truncLen`), respectively  $n=54$  (34); 59 (29); and 14 (5) samples discarded. These three  
787 exceptions suggest that `truncLen`-based trimming can retain substantially more reads and samples for FASTQs with lower  
788 quality reads, which could impact relative abundances (Fig. B1).

789

790 Figures:



791

792 **Figure B1. Relative abundances using different DADA2 trimming methods and the NifMAP OTU pipeline.** *nifH* cluster relative  
793 abundances are shown for each study when processed using the NifMAP OTU pipeline (Angel et al., 2018) or by the *nifH* workflow using  
794 two methods for trimming reads, quality-based (truncQ) or fixed-length (truncLen). ASV or OTU abundances for the samples in a study  
795 were pooled to calculate the relative abundances shown. The three results for each study were calculated using only the samples that were  
796 retained by both runs of the *nifH* workflow. Shiozaki\_2017 and Shiozaki\_2018GBC used mixed-orientation sequencing libraries and could  
797 not be processed by NifMAP.

798

799 Tables:

800 **Table B1. Impact of read trimming method on workflow outputs.** The table compares the *nifH* ASV database, generated using  
801 fixed-length read trimming (truncLen for DADA2 filterAndTrim), to an alternative database for which reads were trimmed at the first  
802 nucleotide where the error rate was >2.5 % (truncQ=16). No other pipeline or post-pipeline parameters were changed.

803

	truncLen	truncQ	% decrease
<b>Samples</b>	944	847	10.4
<b>ASVs</b>	9383	1997	78.7

<b>Reads</b>	43.0E+6	26.3E+6	38.9
--------------	---------	---------	------

804

805

### 806 **Appendix C: Comparison of communities from the workflow to previous studies**

807 Prior to DADA2 (Callhan et al. 2016) and other approaches that distinguish fine-scale variation from sequencing errors  
808 (Eren et al. 2014, Edgar 2016b, Amir et al. 2017), most amplicon studies—for 16S rRNA as well as functional  
809 genes—processed their sequencing data into operational taxonomic units (OTUs). Usually this meant *de novo* clustering the  
810 amplicon sequences at 97 % nucleotide identity and using a representative sequence from each of the OTUs (clusters) for  
811 subsequent analyses. For 16S rRNA genes, it is known that PCR artifacts and sequencing errors can inflate the number of  
812 OTUs and cause diversity to be overestimated (Quince et al., 2009; Eren et al., 2013). For *nifH* amplicon data, these issues  
813 have been mitigated in previously published OTU analyses by analyzing broad diazotroph groups (Table C1).

814

815 To demonstrate whether communities derived from the workflow differ substantially from those previously published, a  
816 comparison was made between the results from the *nifH* workflow and another *nifH* pipeline, NifMAP (Angel et al. 2018).  
817 NifMAP is an OTU pipeline that uses hidden Markov models in an attempt to distinguish true *nifH* sequences from orthologs  
818 often mistaken for *nifH*. NifMAP was used to generate proxies for most of the 21 studies since complete OTU sequences and  
819 abundances were not available for the 19 original studies. Using NifMAP for all studies was more systematic than trying to  
820 reproduce the original results which depended on different software and methods for quality filtering. Additionally, the  
821 workflow and NifMAP both use CART (Frank et al. 2016) to identify *nifH* clusters enabling the cross-comparison of major  
822 *nifH* groups. Both also distinguish *nifH* from orthologs, the workflow using `classifyNifH.sh` described in section 2.3.3). Only  
823 the samples that were processed by both the workflow and NifMAP were compared (n=902).

824

825 The main result was that similar diazotroph communities were detected by the *nifH* workflow and NifMAP (Fig. B1). For  
826 every study they agreed on the two most abundant *nifH* subclusters, usually with  $\leq 3$  % difference between the relative  
827 abundances from the workflow and NifMAP. These results suggest that comparisons between new and previously published  
828 *nifH* amplicon studies are possible, especially if both use similarly broad taxonomic levels, e.g., *nifH* subclusters.

829

830 However, for two studies there were clear differences between the *nifH* workflow and NifMAP that speak to the utility of the  
831 workflow. For Hallstrom\_2022 the workflow detected additional *nifH* subclusters, mainly 3 and 1G, and for Sato\_20201 the  
832 workflow detected 1G and 1A at much higher levels (Fig. B1). These compositional differences likely stemmed from vastly  
833 greater numbers of reads retained by the workflow compared to NifMAP (1034 % and 264 % more reads, respectively for  
834 the two studies; Table C1). The NifMAP logs revealed that poor read quality caused NifMAP to discard the majority of reads  
835 in the first two steps. Only 10% of the Hallstrom\_2022 reads could be merged, the lowest of any study (median 78 %, range  
836 10–94 %), and 56 % of the reads from Sato\_2021. The merged reads were short for both Hallstrom\_2022 (mean 174 nt) and



837 Sato\_2021 (198 nt) in comparison to all studies (median of 307 nt). NifMAP then discarded, respectively, 66 % and 58 % of  
 838 the merged reads due to lengths < 200 nt. In comparison, for Hallstrom\_2022 the workflow discarded most reads during  
 839 DADA2 filterAndTrim (using truncLen) due to sequencing errors but discarded few reads during mergePairs (Fig. 3 and  
 840 Table 4). This suggests that DADA2 denoising worked very well for this dataset because the forward and reverse ASVs were  
 841 allowed at most one mismatch in their overlapping region (Table 2). In contrast, Sato\_2021 had substantial losses of reads  
 842 during both filterAndTrim and mergePairs (Fig. 3 and Table 4). Together these results indicate that the *nifH* workflow can  
 843 potentially retain more reads than NifMAP, particularly when data quality is low, with noticeable impacts on community  
 844 composition.

845

846 Although community compositions from the workflow and NifMAP were mainly similar (Fig. B1), the workflow tended to  
 847 retain more of the sequencing reads (Table C1). For 9 of the 18 studies analyzed by both the workflow and NifMAP, there  
 848 was <10 % difference in the number of reads retained into final sequences (ASVs or OTUs; Table C1). However, 6 of the  
 849 other 9 studies had more reads retained by the workflow (14–1034 %) and 3 had more reads retained by NifMAP (10–23 %).  
 850 Although the workflow retained more reads, usually there were fewer ASVs than OTUs despite compression from clustering  
 851 at 97 % nucleotide identity (Table C1). This is consistent with the known limitations of OTUs mentioned earlier, errors and  
 852 overestimated diversity.

853

854

855 Tables:

856

857 **Table C1.** Summary of the total reads and final sequences obtained by the workflow (ASVs) and NifMAP (OTUs) applied to the same  
 858 samples. A total of 902 of 944 samples in the *nifH* ASV database were compared. This excludes 42 samples from Shiozaki\_2017 and  
 859 Shiozaki\_2018GBC that used mixed-orientation sequencing libraries and could not be processed by NifMAP. The Change (%) column is  
 860 relative to reads in OTUs. OTUs in column 6 count clusters (97 % nucleotide identity). \*: The original publication analyzed OTUs.  
 861

Study ID	Samples compared	Reads (K)			Sequences	
		In OTUs	In ASVs	Change (%)	OTUs	ASVs
AK2HI	43	1319	1259	4.6	987	283
BentzonTilia_2015*	54	220	171	22.6	1043	352
Ding_2021*	32	1358	1446	-6.5	1362	435
Gradoville_2020 (G1,G2)*	162	3200	3304	-3.3	642	333
Hallstrom_2021	82	4531	10,216	-125.5	14,606	6403
Hallstrom_2022*	59	455	5155	-1033.8	91	165
Harding_2018*	88	1384	1579	-14.1	1715	842
Mulholland_2018	28	2527	2439	3.5	1706	549

NEMO	54	1830	1665	9.0	591	177
Raes_2020	131	7668	7793	-1.6	1421	395
Sato_2021	28	106	388	-264.1	141	169
Selden_2021	10	405	445	-9.9	217	60
Shiozaki_2018LNO*	20	618	913	-47.8	929	283
Shiozaki_2020	14	946	1935	-104.7	1664	123
Tang_2020*	6	229	196	14.2	235	35
TurkKubo_2021*	59	2011	1976	1.8	305	74
Wu_2019*	18	801	734	8.3	504	102
Wu_2021*	14	749	674	10.0	1315	180

862

863

#### 864 Author Contributions

865 KTK and MM designed the study with input from SC and MMM. JM created and optimized the DADA2 pipeline for *nifH*  
866 amplicon analyses. JM and MM developed the post-pipeline workflow. MM and JM compiled the database, retrieved  
867 environmental data from CMAP, and analyzed the database. MM, JM and KTK wrote the manuscript with input from  
868 MMM, SC, and JPZ.

#### 869 Competing Interests

870 No competing interest is declared.

#### 871 Acknowledgements

872 We gratefully acknowledge Mohammad Ashkezari and the Simons CMAP team, Stefan Green (Rush University) and the  
873 DNA genomics core at University of Illinois at Chicago, Irina Shilova, Julie Robidart and Grace Reed for NEMO sampling  
874 and sample processing, and Angelicque White (University of Hawaii, Manoa) and Mary R. Gradoville (Columbia River  
875 Inter-Tribal Fish Commission) for AK2HI sampling and sample processing. We would like to thank the authors who directly  
876 provided access to sequences: Takuhei Shiozaki (Shiozaki\_2017 and Shiozaki\_2018GBC) and Jun Sun, Changling Ding, and  
877 Chao Wu (Wu\_2021). This work was supported by grants from the National Science Foundation to KTK (OCE-2023498)  
878 and the Simons Foundation to JPZ (Simons Collaboration on Ocean Processes and Ecology, Award ID 724220).

879

## 880 References

- 881 Amir, A. McDonald, D., Navas-Molina, J. A., Kopylova, E., Morton, J. T., Zech Xu, Z., Kightley, E. P., Thompson, L. R.,  
882 Hyde, E. R., Gonzalez, A., and Knight, R.: Deblur rapidly resolves single-nucleotide community sequence patterns,  
883 *mSystems*, 2, 10.1128/msystems.00191-16, 2017.
- 884 Angel, R., Nepel, M., Panholzl, C., Schmidt, H., Herbold, C. W., Eichorst, S. A., and Wobken, D.: Evaluation of Primers  
885 Targeting the Diazotroph Functional Gene and Development of NifMAP - A Bioinformatics Pipeline for Analyzing nifH  
886 Amplicon Data, *Front Microbiol*, 9, 703, 10.3389/fmicb.2018.00703, 2018.
- 887 Ashkezari, M. D., Hagen, N. R., Denholtz, M., Neang, A., Burns, T. C., Morales, R. L., Lee, C. P., Hill, C. N., and Armbrust,  
888 E. V.: Simons Collaborative Marine Atlas Project (Simons CMAP): An open-source portal to share, visualize, and analyze  
889 ocean data, *Limnol. Oceanogr.: Methods*, 19, 488-496, 2021.
- 890 Benavides, M., Conradt, L., Bonnet, S., Berman-Frank, I., Barrillon, S., Petrenko, A., and Dogliolii, A.: Fine-scale sampling  
891 unveils diazotroph patchiness in the South Pacific Ocean, *ISME Communications*, 1, 3, 2021.
- 892 Bentzon-Tilia, M., Traving, S. J., Mantikci, M., Knudsen-Leerbeck, H., Hansen, J. L. S., Markager, S., and Riemann, L.:  
893 Significant N<sub>2</sub> fixation by heterotrophs, photoheterotrophs and heterocystous cyanobacteria in two temperate estuaries, *ISME*  
894 *J*, 9, 273-285, 2015.
- 895 Blais, M., Tremblay, J. É., Jungblut, A. D., Gagnon, J., Martin, J., Thaler, M., and Lovejoy, C.: Nitrogen fixation and  
896 identification of potential diazotrophs in the Canadian Arctic, *Global Biogeochem. Cy.*, 26, GB3022,  
897 10.1029/2011gb004096, 2012.
- 898 Bostrom, K. H., Riemann, L., Kuhl, M., and Hagstrom, A.: Isolation and gene quantification of heterotrophic N<sub>2</sub>-fixing  
899 bacterioplankton in the Baltic Sea, *Environ. Microbiol.*, 9, 152-164, doi:10.1111/j.1462-2920.2006.01124.x, 2007.
- 900 Cabello, A. M., Turk-Kubo, K. A., Hayashi, K., Jacobs, L., Kudela, R. M., and Zehr, J. P.: Unexpected presence of the  
901 nitrogen-fixing symbiotic cyanobacterium UCYN-A in Monterey Bay, California, *J Phycol*, 56, 1521-1533,  
902 10.1111/jpy.13045, 2020.
- 903 Callahan, B. J., McMurdie, P. J., Rosen, M. J., Han, A. W., Johnson, A. J., and Holmes, S. P.: DADA2: High-resolution  
904 sample inference from Illumina amplicon data, *Nat Methods*, 13, 581-583, 10.1038/nmeth.3869, 2016.
- 905 Callahan, B.J., McMurdie, P. J., and Holmes, S. P.: Exact sequence variants should replace operational taxonomic units in  
906 marker-gene data analysis, *ISME J*, 11, 12, 2639–2643, 2017.
- 907 Capone, D. G., Burns, J. A., Montoya, J. P., Subramaniam, A., Mahaffey, C., Gunderson, T., Michaels, A. F., and Carpenter,  
908 E. J.: Nitrogen fixation by *Trichodesmium* spp.: An important source of new nitrogen to the tropical and subtropical North  
909 Atlantic Ocean, *Global Biogeochem. Cy.*, 19, GB2024: 2021-2017, 2005.
- 910 Carpenter, E. J. and Capone, D. G.: Nitrogen in the marine environment, Academic Press, New York, 900 pp.1983.
- 911 Carpenter, E. J. and Foster, R. A.: Marine symbioses, in: *Cyanobacteria in Symbiosis*, edited by: Rai, A. N., Bergman, B.,  
912 and Rasmussen, U., Kluwer Academic Publishers, The Netherlands, 11-18, 2002.
- 913 Cheung, S., Xia, X., Guo, C., and Liu, H.: Diazotroph community structure in the deep oxygen minimum zone of the Costa  
914 Rica Dome, *J Plankton Res*, 38, 380-391, 2016.

915 Cheung, S., Zehr, J. P., Xia, X., Tsurumoto, C., Endo, H., Nakaoka, S. I., Mak, W., Suzuki, K., and Liu, H.: Gamma4: a  
916 genetically versatile Gammaproteobacterial *nifH* phylotype that is widely distributed in the North Pacific Ocean, Environ  
917 Microbiol, 23, 4246-4259, 10.1111/1462-2920.15604, 2021.

918 Coale, T. H., Loconte, V., Turk-Kubo, K. A., Vanslebrouck, B., Mak, W. K. E., Cheung, S., Ekman, A., Chen, J. H.,  
919 Hagino, K., Takano, Y., Nishimura, T., Adachi, M., Le Gros, M., Larabell, C., and Zehr, J. P.: Nitrogen-fixing organelle in a  
920 marine alga, Science, 384, 217-222, 10.1126/science.adk1075, 2024.

921 Delmont, T. O., Quince, C., Shaiber, A., Esen, Ö. C., Lee, S. T., Rappé, M. S., MacLellan, S. L., Lückner, S., and Eren, A. M.:  
922 Nitrogen-fixing populations of Planctomycetes and Proteobacteria are abundant in surface ocean metagenomes, Nature  
923 microbiology, 3, 804-813, 2018.

924 Delmont, T. O., Karlusich, J. J. P., Veseli, I., Fuessel, J., Eren, A. M., Foster, R. A., Bowler, C., Wincker, P., and Pelletier, E.:  
925 Heterotrophic bacterial diazotrophs are more abundant than their cyanobacterial counterparts in metagenomes covering most  
926 of the sunlit ocean, ISME J, 16, 927-936, 2022.

927 Ding, C., Wu, C., Li, L., Pujari, L., Zhang, G., and Sun, J.: Comparison of Diazotrophic Composition and Distribution in the  
928 South China Sea and the Western Pacific Ocean, Biology (Basel), 10, 10.3390/biology10060555, 2021.

929 Edgar, R.: UCHIME2: improved chimera prediction for amplicon sequencing, BioRxiv, doi.org/10.1101/074252, 2016a.

930 Edgar, R.: UNOISE2: improved error-correction for Illumina 16S and ITS amplicon sequencing, BioRxiv,  
931 doi.org/10.1101/081257, 2016b.

932 Eren, A. M., Vineis, J. H., Morrison, H. G., and Sogin, M. L.: A filtering method to generate high quality short reads using  
933 Illumina paired-end technology, PLOS ONE, 8, 6, e66643, 10.1371/journal.pone.0066643, 2013.

934 Eren, A. M., Morrison, H. G., Lescault, P. J., Reveillaud, J., Vineis, J. H. and Sogin, M. L.: Minimum entropy  
935 decomposition: Unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences, ISME J, 9, 4,  
936 968-979, doi.org/10.1038/ismej.2014.195, 2014.

937 Falcón, L., Cipriano, F., Chistoserdov, A., and Carpenter, E.: Diversity of diazotrophic unicellular cyanobacteria in the  
938 tropical North Atlantic Ocean, Appl Environ Microbiol, 68, 5760, 2002.

939 Falcón, L., Carpenter, E., Cipriano, F., Bergman, B., and Capone, D.: N<sub>2</sub> fixation by unicellular bacterioplankton from the  
940 Atlantic and Pacific Oceans: phylogeny and in situ rates, Appl Environ Microbiol, 70, 765-770, 2004.

941 Farnelid, H., Oberg, T., and Riemann, L.: Identity and dynamics of putative N<sub>2</sub>-fixing picoplankton in the Baltic Sea proper  
942 suggest complex patterns of regulation, Environmental Microbiology Reports, 1, 145-154,  
943 10.1111/j.1758-2229.2009.00021.x, 2009.

944 Farnelid, H., Andersson, A. F., Bertilsson, S., Al-Soud, W. A., Hansen, L. H., Sørensen, S., Steward, G. F., Hagström, A.,  
945 and Riemann, L.: Nitrogenase gene amplicons from global marine surface waters are dominated by genes of  
946 non-cyanobacteria, PLOS ONE, 6, e19223, 10.1371/journal.pone.0019223, 2011.

947 Fernandez, C., Farias, L., and Ulloa, O.: Nitrogen fixation in denitrified marine waters, PLOS ONE, 6, e20539,  
948 10.1371/journal.pone.0020539, 2011.

949 Fernandez-Mendez, M., Turk-Kubo, K. A., Buttigieg, P. L., Rapp, J. Z., Krumpen, T., Zehr, J. P., and Boetius, A.: Diazotroph  
950 Diversity in the Sea Ice, Melt Ponds, and Surface Waters of the Eurasian Basin of the Central Arctic Ocean, *Front Microbiol*,  
951 7, 1-18, 10.3389/fmicb.2016.01884, 2016.

952 Frank, I. E., Turk-Kubo, K. A., and Zehr, J. P.: Rapid annotation of *nifH* gene sequences using classification and regression  
953 trees facilitates environmental functional gene analysis, *Env Microbiol Rep*, 8, 905-916, 2016.

954 Gaby, J. C. and Buckley, D. H.: A global census of nitrogenase diversity, *Environ Microbiol*, 13, 1790-1799,  
955 10.1111/j.1462-2920.2011.02488.x, 2011.

956 Goto, M., Ando, S., Hachisuka, Y., and Yoneyama, T.: Contamination of diverse *nifH* and *nifH*-like DNA into commercial  
957 PCR primers, *FEMS Microbiol Lett*, 246, 33-38, 10.1016/j.femsle.2005.03.042, 2005.

958 Gradoville, M. R., Bombar, D., Crump, B. C., Letelier, R. M., Zehr, J. P., and White, A. E.: Diversity and activity of  
959 nitrogen-fixing communities across ocean basins, *Limnol Oceanogr*, 62, 1895-1909, 2017.

960 Gradoville, M. R., Farnelid, H., White, A. E., Turk-Kubo, K. A., Stewart, B., Ribalet, F., Ferrón, S., Pinedo-Gonzalez, P.,  
961 Armbrust, E. V., Karl, D. M., John, S., and Zehr, J. P.: Latitudinal constraints on the abundance and activity of the  
962 cyanobacterium UCYN-A and other marine diazotrophs in the North Pacific, *Limnol Oceanogr*, 65, 1858-1875,  
963 10.1002/lno.11423, 2020.

964 Green, S. J., Venkatramanan, R., and Naqib, A.: Deconstructing the polymerase chain reaction: Understanding and correcting  
965 bias associated with primer degeneracies and primer-template mismatches, *PLOS ONE*, 10, e0128122,  
966 doi:10.1371/journal.pone.0128122, 2015.

967 Hallstrøm, S., Benavides, M., Salamon, E. R., Arístegui, J., and Riemann, L.: Activity and distribution of diazotrophic  
968 communities across the Cape Verde Frontal Zone in the Northeast Atlantic Ocean, *Biogeochem*, 1-19, 2022a.

969 Hallstrøm, S., Benavides, M., Salamon, E. R., Evans, C. W., Potts, L. J., Granger, J., Tobias, C. R., Moisander, P. H., and  
970 Riemann, L.: Pelagic N<sub>2</sub> fixation dominated by sediment diazotrophic communities in a shallow temperate estuary, *Limnol*  
971 *Oceanogr*, 67, 364-378, 2022b.

972 Halm, H., Lam, P., Ferdelman, T. G., Lavik, G., Dittmar, T., LaRoche, J., D'Hondt, S., and Kuypers, M. M.: Heterotrophic  
973 organisms dominate nitrogen fixation in the South Pacific Gyre, *ISME J*, 6, 1238-1249, 10.1038/ismej.2011.182, 2012.

974 Harding, K., Turk-Kubo, K. A., Sipler, R. E., Mills, M. M., Bronk, D. A., and Zehr, J. P.: Symbiotic unicellular  
975 cyanobacteria fix nitrogen in the Arctic Ocean, *Proc Natl Acad Sci U S A*, 115, 13371-13375, 10.1073/pnas.1813658115,  
976 2018.

977 Heller, P., Tripp, H. J., Turk-Kubo, K., and Zehr, J. P.: ARBitrator: a software pipeline for on-demand retrieval of  
978 auto-curated *nifH* sequences from GenBank, *Bioinformatics*, 10.1093/bioinformatics/btu417, 2014.

979 Jickells, T., Buitenhuis, E., Altieri, K., Baker, A., Capone, D., Duce, R., Dentener, F., Fennel, K., Kanakidou, M., and  
980 LaRoche, J.: A reevaluation of the magnitude and impacts of anthropogenic atmospheric nitrogen inputs on the ocean,  
981 *Global Biogeochem. Cy.*, 31, 289-305, 2017.

982 Langlois, R. J., LaRoche, J., and Raab, P. A.: Diazotrophic diversity and distribution in the tropical and subtropical Atlantic  
983 Ocean, *Appl Environ Microbiol*, 71, 7910-7919, 10.1128/AEM.71.12.7910-7919.2005, 2005.

- 984 Langlois, R., Großkopf, T., Mills, M., Takeda, S., and LaRoche, J.: Widespread distribution and expression of Gamma A  
985 (UMB), an uncultured, diazotrophic,  $\gamma$ -proteobacterial nifH phylotype, PLOS ONE, 10, e0128912, 2015.
- 986 Liu, J., Zhou, L., Li, J., Lin, Y., Ke, Z., Zhao, C., Liu, H., Jiang, X., He, Y., and Tan, Y.: Effect of mesoscale eddies on  
987 diazotroph community structure and nitrogen fixation rates in the South China Sea, Regional Studies in Marine Science, 35,  
988 101106, 2020.
- 989 Löscher, C. R., Großkopf, T., Desai, F. D., Gill, D., Schunck, H., Croot, P. L., Schlosser, C., Neulinger, S. C., Pinnow, N.,  
990 and Lavik, G.: Facets of diazotrophy in the oxygen minimum zone waters off Peru, ISME J, 8, 2180-2192, 2014.
- 991 Luo, Y. W., Doney, S. C., Anderson, L. A., Benavides, M., Berman-Frank, I., Bode, A., Bonnet, S., Boström, K. H., Böttjer,  
992 D., Capone, D. G., Carpenter, E. J., Chen, Y. L., Church, M. J., Dore, J. E., Falcón, L. I., Fernández, A., Foster, R. A.,  
993 Furuya, K., Gómez, F., Gundersen, K., Hynes, A. M., Karl, D. M., Kitajima, S., Langlois, R. J., LaRoche, J., Letelier, R. M.,  
994 Marañón, E., McGillicuddy, D. J., Moisaner, P. H., Moore, C. M., Mouriño-Carballido, B., Mulholland, M. R., Needoba, J.  
995 A., Orcutt, K. M., Poulton, A. J., Rahav, E., Raimbault, P., Rees, A. P., Riemann, L., Shiozaki, T., Subramaniam, A., Tyrrell,  
996 T., Turk-Kubo, K. A., Varela, M., Villareal, T. A., Webb, E. A., White, A. E., Wu, J., and Zehr, J. P.: Database of diazotrophs  
997 in global ocean: abundance, biomass and nitrogen fixation rates, Earth System Science Data, 4, 47-73,  
998 10.5194/essd-4-47-2012, 2012.
- 999 Manley, L. J., Ma, D., and Levine, S. S.: Monitoring error rates In Illumina sequencing, J Biomol Tech, 27, 4, 125-128,  
1000 10.7171/jbt.16-2704-002, 2016.
- 1001 Martin, M.: Cutadapt removes adapter sequences from high-throughput sequencing reads, EMBnet, 17, 10-12, 2011.
- 1002 Messer, L. F., Mahaffey, C., Robinson, C. M., Jeffries, T. C., Baker, K. G., Isaksson, J. B., Ostrowski, M., Doblin, M. A.,  
1003 Brown, M. V., and Seymour, J. R.: High levels of heterogeneity in diazotroph diversity and activity within a putative hotspot  
1004 for marine nitrogen fixation, ISME J, 1499-1513, 2015.
- 1005 Moisaner, P. H., Beinart, R. A., Voss, M., and Zehr, J. P.: Diversity and abundance of diazotrophic microorganisms in the  
1006 South China Sea during intermonsoon, ISME J, 2, 954-967, 10.1038/ismej.2008.51, 2008.
- 1007 Moisaner, P. H., Serros, T., Paerl, R. W., Beinart, R. A., and Zehr, J. P.: Gammaproteobacterial diazotrophs and nifH gene  
1008 expression in surface waters of the South Pacific Ocean, ISME J, 8, 1962-1973, 10.1038/ismej.2014.49, 2014.
- 1009 Moisaner, P. H., Benavides, M., Bonnet, S., Berman-Frank, I., White, A. E., and Riemann, L.: Chasing after  
1010 non-cyanobacterial nitrogen fixation in marine pelagic environments, Front Microbiol, 8, 1736, 2017.
- 1011 Moonsamy, P. V., Williams, T., Bonella, P., Holcomb, C. L., Høglund, B. N., Hillman, G., Goodridge, D., Turenchalk, G. S.,  
1012 Blake, L. A., Daigle, D. A., Simen, B. B., Hamilton, A., May, A. P., and Erlich, H. A.: High throughput HLA genotyping  
1013 using 454 sequencing and the Fluidigm Access Array System for simplified amplicon library preparation, Tissue Antigens,  
1014 81, 141-149, 10.1111/tan.12071, 2013.
- 1015 Morando, M., Magasin, J., Cheung, S., Mills, M. M., Zehr, J. P., and Turk-Kubo, K. A.: nifH ASV database in Global  
1016 biogeography of  $N_2$ -fixing microbes: nifH amplicon database and analytics workflow, Figshare [dataset],  
1017 <https://doi.org/10.6084/m9.figshare.23795943.v2>, 2024a.
- 1018 Morando, M., Magasin, J., Cheung, S., Mills, M. M., Zehr, J. P., and Turk-Kubo, K. A.: DADA2 nifH pipeline in Global  
1019 biogeography of  $N_2$ -fixing microbes: nifH amplicon database and analytics workflow, GitHub [code],  
1020 [https://github.com/jdmagasin/nifH\\_amplicons\\_DADA2](https://github.com/jdmagasin/nifH_amplicons_DADA2), 2024b.

- 1021 Morando, M., Magasin, J., Cheung, S., Mills, M. M., Zehr, J. P., and Turk-Kubo, K. A.: nifH ASV workflow in Global  
1022 biogeography of N<sub>2</sub>-fixing microbes: nifH amplicon database and analytics workflow, GitHub [code],  
1023 <https://github.com/jdmagasin/nifH-ASV-workflow>, 2024c.
- 1024 Mulholland, M. R., Bernhardt, P. W., Widner, B. N., Selden, C. R., Chappell, P. D., Clayton, S., Mannino, A., and Hyde, K.:  
1025 High Rates of N<sub>2</sub> Fixation in Temperate, Western North Atlantic Coastal Waters Expand the Realm of Marine Diazotrophy,  
1026 *Global Biogeochem. Cy.*, 33, 826-840, 10.1029/2018gb006130, 2019.
- 1027 Pierella Karlusich, J. J., Pelletier, E., Lombard, F., Carsique, M., Dvorak, E., Colin, S., Picheral, M., Cornejo-Castillo, F. M.,  
1028 Acinas, S. G., Pepperkok, R., Karsenti, E., de Vargas, C., Wincker, P., Bowler, C., and Foster, R. A.: Global distribution  
1029 patterns of marine nitrogen-fixers by imaging and molecular methods, *Nat Commun*, 12, 1-18, 10.1038/s41467-021-24299-y,  
1030 2021.
- 1031 Quince, C., Lanzén, A., Curtis, T. P., Davenport, R. J., Hall, N., Head, I. M., Read, L. F., and Sloan, W. T.: Accurate  
1032 determination of microbial diversity from 454 pyrosequencing data, *Nat Methods*, 6, 639–641, doi.org/10.1038/nmeth.1361,  
1033 2009.
- 1034 Raes, E. J., Van de Kamp, J., Bodrossy, L., Fong, A. A., Riekenberg, J., Holmes, B. H., Erler, D. V., Eyre, B. D., Weil, S. S.,  
1035 and Waite, A. M.: N<sub>2</sub> fixation and new insights into nitrification from the ice-edge to the equator in the South Pacific Ocean,  
1036 *Frontiers in Marine Science*, 7, 1-20, 2020.
- 1037 Rho, M., Tang, H., and Ye, Y.: FragGeneScan: predicting genes in short and error-prone reads, *Nucleic Acids Res*, 38, e191,  
1038 10.1093/nar/gkq747, 2010.
- 1039 Rognes, T., Flouri, T., Nichols, B., Quince, C., and Mahe, F.: VSEARCH: a versatile open source tool for metagenomics,  
1040 *PeerJ*, 4, e2584, 10.7717/peerj.2584, 2016.
- 1041 Sato, T., Shiozaki, T., Taniuchi, Y., Kasai, H., and Takahashi, K.: Nitrogen Fixation and Diazotroph Community in the  
1042 Subarctic Sea of Japan and Sea of Okhotsk, *Journal of Geophysical Research: Oceans*, 126, e2020JC017071, 2021.
- 1043 Schlessman, J. L., Woo, D., Joshua-Tor, L., Howard, J. B., and Rees, D. C.: Conformational variability in structures of the  
1044 nitrogenase iron proteins from *Azotobacter vinelandii* and *Clostridium pasteurianum*, *J Mol Biol*, 280, 4, 669-685, 1998.
- 1045 Selden, C. R., Chappell, P. D., Clayton, S., Macías-Tapia, A., Bernhardt, P. W., and Mulholland, M. R.: A coastal N<sub>2</sub> fixation  
1046 hotspot at the Cape Hatteras front: Elucidating spatial heterogeneity in diazotroph activity via supervised machine learning,  
1047 *Limnol Oceanogr*, 66, 1832-1849, 2021.
- 1048 Shao, Z. and Luo, Y. W.: Controlling factors on the global distribution of a representative marine heterotrophic diazotroph  
1049 phylotype (Gamma A), *Biogeosciences*, 19, 2939-2952, 2022.
- 1050 Shao, Z., Xu, Y., Wang, H., Luo, W., Wang, L., Huang, Y., Agawin, N. S. R., Ahmed, A., Benavides, M., Bentzon-Tilia, M.,  
1051 and Berman-Frank, I.: Global Oceanic Diazotroph Database Version 2 and Elevated Estimate of Global N<sub>2</sub> Fixation, *Earth*  
1052 *System Science Data*, 15, 2023.
- 1053 Shilova, I., Mills, M., Robidart, J., Turk-Kubo, K., Björkman, K., Kolber, Z., Rapp, I., van Dijken, G., Church, M., and  
1054 Arrigo, K.: Differential effects of nitrate, ammonium, and urea as N sources for microbial communities in the North Pacific  
1055 Ocean, *Limnol Oceanogr*, 62, 2550-2574, 2017.
- 1056 Shiozaki, T., Fujiwara, A., Inomura, K., Hirose, Y., Hashihama, F., and Harada, N.: Biological nitrogen fixation detected  
1057 under Antarctic sea ice, *Nature Geoscience*, 13, 729–732, 2020.



- 1058 Shiozaki, T., Bombar, D., Riemann, L., Hashihama, F., Takeda, S., Yamaguchi, T., Ehama, M., Hamasaki, K., and Furuya,  
1059 K.: Basin scale variability of active diazotrophs and nitrogen fixation in the North Pacific, from the tropics to the subarctic  
1060 Bering Sea, *Global Biogeochem. Cy.*, 31, 996-1009, 10.1002/2017gb005681, 2017.
- 1061 Shiozaki, T., Fujiwara, A., Ijichi, M., Harada, N., Nishino, S., Nishi, S., Nagata, T., and Hamasaki, K.: Diazotroph  
1062 community structure and the role of nitrogen fixation in the nitrogen cycle in the Chukchi Sea (western Arctic Ocean),  
1063 *Limnol Oceanogr*, 63, 2191-2205, 10.1002/lno.10933, 2018a.
- 1064 Shiozaki, T., Bombar, D., Riemann, L., Sato, M., Hashihama, F., Kodama, T., Tanita, I., Takeda, S., Saito, H., Hamasaki, K.,  
1065 and Furuya, K.: Linkage between dinitrogen fixation and primary production in the oligotrophic South Pacific Ocean, *Global*  
1066 *Biogeochem. Cy.*, 32, 1028-1044, 2018b.
- 1067 Tang, W., Li, Z., and Cassar, N.: Machine learning estimates of global marine nitrogen fixation, *Journal of Geophysical*  
1068 *Research: Biogeosciences*, 124, 717-730, 2019.
- 1069 Tang, W., Cerdan-Garcia, E., Berthelot, H., Polyviou, D., Wang, S., Baylay, A., Whitby, H., Planquette, H., Mowlem, M.,  
1070 Robidart, J., and Cassar, N.: New insights into the distributions of nitrogen fixation and diazotrophs revealed by  
1071 high-resolution sensing and sampling methods, *ISME J*, 14, 2514-2526, 10.1038/s41396-020-0703-6, 2020.
- 1072 Taylor, L. J., Abbas, A., and Bushman, F. D.: grabseqs: Simple downloading of reads and metadata from multiple  
1073 next-generation sequencing data repositories, *Bioinformatics*, doi.org/10.1093/bioinformatics/btaa167, 2020.
- 1074 Turk, K., Rees, A. P., Zehr, J. P., Pereira, N., Swift, P., Shelley, R., Lohan, M., Woodward, E. M. S., and Gilbert, J.: Nitrogen  
1075 fixation and nitrogenase (nifH) expression in tropical waters of the eastern North Atlantic, *ISME J*, 5, 1201-1212,  
1076 10.1038/ismej.2010.205, 2011.
- 1077 Turk-Kubo, K. A., Karamchandani, M., Capone, D. G., and Zehr, J. P.: The paradox of marine heterotrophic nitrogen  
1078 fixation: abundances of heterotrophic diazotrophs do not account for nitrogen fixation rates in the Eastern Tropical South  
1079 Pacific, *Environ Microbiol*, 16, 3095-3114, 10.1111/1462-2920.12346, 2014.
- 1080 Turk-Kubo, K. A., Farnelid, H. M., Shilova, I. N., Henke, B., and Zehr, J. P.: Distinct ecological niches of marine symbiotic  
1081 N<sub>2</sub>-fixing cyanobacterium *Candidatus Atelocyanobacterium thalassa* sublineages, *J Phycol*, 53, 451-461, 10.1111/jpy.12505,  
1082 2017.
- 1083 Turk-Kubo, K. A., Mills, M. M., Arrigo, K. R., van Dijken, G., Henke, B. A., Stewart, B., Wilson, S. T., and Zehr, J. P.:  
1084 UCYN-A/haptophyte symbioses dominate N<sub>2</sub> fixation in the Southern California Current System, *ISME Communications*, 1,  
1085 1-13, 2021.
- 1086 Turk-Kubo, K. A., Gradoville, M. R., Cheung, S., Cornejo-Castillo, F., Harding, K. J., Morando, M., Mills, M., and Zehr, J.  
1087 P.: Non-cyanobacterial diazotrophs: Global diversity, distribution, ecophysiology, and activity in marine waters, *FEMS*  
1088 *Microbiol Rev*, 10.1093/femsre/fuac046, 2022.
- 1089 Villareal, T. A.: Widespread occurrence of the *Hemiaulus*-cyanobacterial symbiosis in the southwest North-Atlantic Ocean,  
1090 *Bulletin of Marine Science*, 54, 1-7, 1994.
- 1091 Wu, C., Kan, J., Liu, H., Pujari, L., Guo, C., Wang, X., and Sun, J.: Heterotrophic Bacteria Dominate the Diazotrophic  
1092 Community in the Eastern Indian Ocean (EIO) during Pre-Southwest Monsoon, *Microb Ecol*, 78, 804-819,  
1093 10.1007/s00248-019-01355-1, 2019.



- 1094 Wu, C., Sun, J., Liu, H., Xu, W., Zhang, G., Lu, H., and Guo, Y.: Evidence of the Significant Contribution of Heterotrophic  
1095 Diazotrophs to Nitrogen Fixation in the Eastern Indian Ocean During Pre-Southwest Monsoon Period, *Ecosyst*, 25,  
1096 1066-1083, 2021.
- 1097 Zani, S.: Application of a nested reverse transcriptase polymerase chain reaction assay for the detection of *nifH* expression in  
1098 Lake George, New York, M. S. Thesis, Rensselaer Polytechnic Institute, 1999.
- 1099 Zehr, J. P. and Capone, D. G.: Changing perspectives in marine nitrogen fixation, *Science*, 368, eaay9514,  
1100 10.1126/science.aay9514, 2020.
- 1101 Zehr, J. and McReynolds, L.: Use of degenerate oligonucleotides for amplification of the *nifH* gene from the marine  
1102 cyanobacterium *Trichodesmium thiebautii*, *Appl Environ Microbiol*, 55, 2522-2526, 1989.
- 1103 Zehr, J., Mellon, M., and Zani, S.: New nitrogen-fixing microorganisms detected in oligotrophic oceans by amplification of  
1104 nitrogenase (*nifH*) genes, *Appl. Environ. Microbiol*, 64, 3444-3450, 1998.
- 1105 Zehr, J. P., Waterbury, J. B., Turner, P. J., Montoya, J. P., Omoregie, E., Steward, G. F., Hansen, A., and Karl, D. M.:  
1106 Unicellular cyanobacteria fix N<sub>2</sub> in the subtropical North Pacific Ocean, *Nature*, 412, 635-638, 2001.  
1107
- 1108 Zehr, J. P., Crumbliss, L. L., Church, M. J., Omoregie, E. O., and Jenkins, B. D.: Nitrogenase genes in PCR and RT-PCR  
1109 reagents: implications for studies of diversity of functional genes, *Biotechniques*, 35, 996-1005, 2003.