

1 **Results from a Multi-Laboratory Ocean Metaproteomic Intercomparison:**

2 **Effects of LC-MS Acquisition and Data Analysis Procedures**

3 ***Participants of the Ocean Metaproteome Intercomparison Consortium:***

4 Mak A. Saito<sup>\*1</sup>, Jaclyn K. Saunders<sup>1\*</sup>, Matthew R. McIlvin<sup>1</sup>, Erin M. Bertrand<sup>2</sup>, John A. Breier<sup>3</sup>,  
5 Margaret Mars Brisbin<sup>1</sup>, Sophie M. Colston<sup>4</sup>, Jaimee R. Compton<sup>4</sup>, Tim J. Griffin<sup>5</sup>, W. Judson  
6 Hervey<sup>4</sup>, Robert L. Hettich<sup>6</sup>, Pratik D. Jagtap<sup>5</sup>, Michael Janech<sup>7</sup>, Rod Johnson<sup>8</sup>, Rick Keil<sup>9</sup>, Hugo  
7 Kleikamp<sup>10</sup>, Dagmar Leary<sup>4</sup>, Lennart Martens<sup>17,18</sup>, J. Scott P. McCain<sup>2,11</sup>, Eli Moore<sup>12</sup>, Subina  
8 Mehta<sup>5</sup>, Dawn M. Moran<sup>1</sup>, Jacqui Neibauer<sup>7</sup>, Benjamin A. Neely<sup>13</sup>, Michael V. Jakuba<sup>1</sup>, Jim  
9 Johnson<sup>5</sup>, Megan Duffy<sup>7</sup>, Gerhard J. Herndl<sup>14</sup>, Richard Giannone<sup>6</sup>, Ryan Mueller<sup>15</sup>, Brook L.  
10 Nunn<sup>9</sup>, Martin Pabst<sup>9</sup>, Samantha Peters<sup>6</sup>, Andrew Rajczewski<sup>5</sup>, Elden Rowland<sup>2</sup>, Brian  
11 Searle<sup>16</sup>, Tim Van Den Bossche<sup>17,18</sup>, Gary J. Vora<sup>4</sup>, Jacob R. Waldbauer<sup>19</sup>, Haiyan Zheng<sup>20</sup>,  
12 Zihao Zhao<sup>14</sup>

13  
14 <sup>1</sup>Woods Hole Oceanographic Institution, Woods Hole, MA, USA

15 <sup>2</sup>Department of Biology, Dalhousie University, Halifax, NS, Canada

16 <sup>3</sup>The University of Texas Rio Grande Valley, Edinburg, TX

17 <sup>4</sup>Center for Bio/Molecular Science & Engineering, Naval Research Laboratory, Washington, DC, USA

18 <sup>5</sup>University of Minnesota at Minneapolis, Minneapolis, Minnesota, USA

19 <sup>6</sup>Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA

20 <sup>7</sup>College of Charleston, Charleston, South Carolina, USA

21 <sup>8</sup>Bermuda Institute of Ocean Sciences, Bermuda

22 <sup>9</sup>University of Washington, Seattle, Washington, USA

23 <sup>10</sup>Department of Biotechnology, Delft University of Technology, Netherlands

24 <sup>11</sup>Department of Biology, Massachusetts Institute of Technology, Cambridge, MA, USA

25 <sup>12</sup>United States Geological Survey, USA

26 <sup>13</sup>National Institute of Standards and Technology, Charleston, South Carolina, USA

27 <sup>14</sup>University of Vienna, Dept. of Functional and Evolutionary Ecology, Austria

28 <sup>15</sup>Oregon State University, Corvallis, Oregon, USA

29 <sup>16</sup>Ohio State University, Columbus, Ohio, USA

30 <sup>17</sup>Department of Biomolecular Medicine, Faculty of Medicine and Health Sciences, Ghent University, 9000  
31 Ghent, Belgium

32 <sup>18</sup>VIB – UGent Center for Medical Biotechnology, VIB, 9000 Ghent, Belgium

33 <sup>19</sup>Department of Geophysical Sciences, University of Chicago, Chicago, Illinois, USA

34 <sup>20</sup>Rutgers University, Piscataway, New Jersey, USA

35 <sup>\*</sup>Present address: University of Georgia, Department of Marine Sciences

36 <sup>\*</sup> corresponding author, msaito@whoi.edu

37  
38  
39 ~~For~~ Submission to Biogeosciences 12/27/2023

40 Revision for Biogeosciences 5/2024

41

42 **Abstract**

43 Metaproteomics is an increasingly popular methodology that provides information regarding the  
44 metabolic functions of specific microbial taxa and has potential for contributing to ocean ecology  
45 and biogeochemical studies. A blinded multi-laboratory intercomparison was conducted to  
46 assess comparability and reproducibility of taxonomic and functional results and their sensitivity  
47 to methodological variables. Euphotic zone samples from the Bermuda Atlantic Time-Series  
48 Study in the North Atlantic Ocean collected by *in situ* pumps and the AUV *Clio* were distributed  
49 with a paired metagenome, and one-dimensional liquid chromatographic data dependent  
50 acquisition mass spectrometry analyses was stipulated. Analysis of mass spectra from seven  
51 laboratories through a common [bioinformatic](#) pipeline identified a shared set of 1056 proteins  
52 from 1395 shared peptides constituents. Quantitative analyses showed good reproducibility:  
53 pairwise regressions of spectral counts between laboratories yielded R<sup>2</sup> values [averaged 0.62](#)  
54 [+/- 0.11 ranging from 0.43 to 0.83](#), and a Sørensen similarity analysis of the top 1,000 proteins  
55 revealed 70-80% similarity between laboratory groups. Taxonomic and functional assignments  
56 showed good coherence between technical replicates and different laboratories. A [bioinformatic](#)  
57 intercomparison study, involving 10 laboratories using 8 software packages  
58 successfully identified thousands of peptides within the complex metaproteomic datasets,  
59 demonstrating the utility of these software tools for ocean metaproteomic research. [Lessons](#)  
60 [learned and potential improvements in methods were described.](#) Future efforts could examine  
61 reproducibility in deeper metaproteomes, examine accuracy in targeted absolute quantitation  
62 analyses, and develop standards for data output formats to improve data interoperability.  
63 Together, these results demonstrate the reproducibility of metaproteomic analyses and their  
64 suitability for microbial oceanography research including integration into global scale ocean  
65 surveys and ocean biogeochemical models.

66

67 **1. Introduction**

68 Microorganisms within the oceans are major contributors to global biogeochemical cycles,  
69 influencing the cycling of carbon, nitrogen, phosphorus, sulfur, iron, cobalt and other elements  
70 (Falkowski et al., 2008; Moran et al., 2022; Worden et al., 2015). 'Omic methodologies can  
71 provide an expansive window into these communities, with genomic approaches characterizing  
72 the diversity and potential metabolisms, and transcriptomic and proteomic methods providing  
73 insights into expression and function of that potential. ~~Similar to other 'omics approaches,~~  
74 ~~proteomics is increasingly being applied to natural ocean environments and the diverse~~  
75 ~~microbial communities within them. When proteomics is applied to such mixed communities, it is~~  
76 ~~generally referred to as metaproteomics (Wilmes and Bond, 2006). Of these, proteomics is~~  
77 ~~increasingly being applied to natural ocean environments—when applied to complex~~  
78 ~~communities with diverse taxa present, the technique is commonly referred to as~~  
79 ~~metaproteomics (Wilmes and Bond, 2006).~~ Metaproteomic samples contain an extraordinary  
80 level of complexity relative to single organism proteomes (at least 1-2 orders of magnitude) due  
81 to the simultaneous presence of many different organisms in widely varying abundances  
82 (McCain and Bertrand, 2019). In particular, ocean metaproteome samples are significantly more  
83 complex than the human proteome, the latter of which is itself considered to be a highly  
84 complex sample (Saito et al., 2019). Proteomics (including metaproteomics) provides a  
85 perspective distinct from other 'omics methods: as a direct measurement of cellular functions it  
86 can be used to examine the diversity of ecosystem biogeochemical capabilities, to determine  
87 the extent of specific nutrient stressors by measurement of transporters or regulatory systems,  
88 to determine cellular resource allocation strategies in-situ, estimate biomass contributions from  
89 specific microbial groups, and even to estimate potential enzyme activity (Bender et al., 2018;  
90 Bergauer et al., 2018; Cohen et al., 2021; Fuchsman et al., 2019; Georges et al., 2014; Hawley  
91 et al., 2014; Held et al., 2021; Leary et al., 2014; McCain et al., 2022; Mikan et al., 2020; Moore

Formatted: Font: Not Italic

Formatted: Font: Not Italic

Formatted: Font: Not Italic

Formatted: Font: Not Italic

Formatted: Font: Not Italic

Formatted: Font: Not Italic

Formatted: Font: Not Italic

92 et al., 2012; Morris et al., 2010; Saito et al., 2020; Sowell et al., 2009; Williams et al., 2012). The  
93 functional perspective that metaproteomics allows is often complementary to metagenomic and  
94 metatranscriptomic analyses and can provide biological insights that are distinct from organisms  
95 studied in the laboratory (Kleiner et al., 2019). Moreover, the measurement of microbial proteins  
96 in environmental samples has improved greatly in recent years, due to the advancements in  
97 nanospray-liquid chromatography and high-resolution mass spectrometry approaches (Mueller  
98 and Pan, 2013; Ram et al., 2005; McIlvin and Saito, 2021).

99 With increasing interest in the measurement of proteins and their biogeochemical  
100 functions within the oceans, the metaproteomic datatype is beginning to establish itself as a  
101 valuable research and monitoring tool. However, given rapid changes in technology and  
102 methods, as well as the overall youth of the metaproteomic field, demonstrating the  
103 reproducibility and robustness of metaproteomic measurements to microbial ecology and  
104 oceanographic communities is an important goal. This is particularly true as applications for  
105 metaproteomics expand in research and monitoring of the changing ocean environment, for  
106 example in global scale efforts such as the developing BioGeoSCAPES program  
107 ([www.biogeoscapes.org](http://www.biogeoscapes.org); (Tagliabue, 2023)), which aims to characterize the ocean metabolism  
108 and nutrient cycles on a changing planet. As a result, there is a pressing need to assess inter-  
109 laboratory consistency, and to understand the impacts of sampling, extraction, mass  
110 spectrometry, and bioinformatic analyses on the biological inferences that can be drawn from  
111 the data.

112 There have been efforts to conduct intercomparisons of metaproteomic analyses in both  
113 biomedical and environmental sample types in recent years that provide precedent for this  
114 study. A recent community best practice effort in ocean metaproteomics data-sharing also  
115 identified major challenges in ocean metaproteomics research, including sampling, extraction,  
116 sample analysis, bioinformatics pipelines, and data sharing, and conducted a quantitative

117 assessment of sample complexity in ocean metaproteome samples (Saito et al., 2019). A  
118 previous benchmark study, driven by the Metaproteomics Initiative (Van Den Bossche et al.,  
119 2021), was the “Critical Assessment of Metaproteome Investigation study” (CAMPI) that  
120 employed a laboratory-assembled microbiome and human fecal microbiome sample to  
121 successfully demonstrate reproducibility of results between laboratories. CAMPI found  
122 robustness in results across datasets, while also observing variability in peptide identifications  
123 largely attributed to sample preparation. This observation was consistent with prior findings on  
124 single organism samples that determined >70% of the variability was due to sample processing,  
125 rather than chromatography and mass spectrometry (Piehowski et al., 2013). Finally, the  
126 Proteomics Informatics Group (iPRG) from the Association of Biomolecular Resources Facilities  
127 (ABRF) conducted a study examining the influence of informatics pipelines on metaproteomics  
128 analyses that found consistency among research groups in taxonomic attributions (Jagtap et al.,  
129 2023), and previous research has demonstrated the impact of database choices on final  
130 functional annotations and biological implications (Timmins-Schiffman et al., 2017).

131 Here we describe the results from the first ocean metaproteomic intercomparison. In this  
132 study, environmental ocean samples were collected from the euphotic zone of the North Atlantic  
133 Ocean and partitioned into subsamples and distributed to an international group of laboratories  
134 (Fig. 1). The study was designed to examine inter-laboratory consistency rather than maximal  
135 capabilities, stipulating one-dimensional chromatographic analyses from each laboratory (with  
136 optional deeper analysis). Users were invited to use their preferred extraction, analytical, and  
137 bioinformatic procedures. The effort focused on the data dependent analysis (DDA) methods,  
138 also known as global proteomics where the targets are unknown and hence there is a discovery  
139 element to the approach. DDA is ~~that are~~ currently common in ocean and other environmental  
140 and biomedical metaproteomics, and ~~its associated~~ spectral abundance units of relative  
141 quantitation, ~~which~~ have been shown to be reproducible in metaproteomics (Kleiner et al., 2017;

142 Pietilä et al., 2022). Blinded results were submitted, compared and discussed at a virtual  
143 community workshop in September of 2021. An additional [bioinformatic](#) pipeline comparison  
144 study was also conducted where participants were provided metaproteomic raw data and  
145 associated metagenomic sequence database files and were encouraged to use the  
146 [bioinformatic](#) pipeline of their choice.

## 147 **2. Methods**

### 148 *2.1 Sample Collection and Metadata*

149 Ocean metaproteome filter samples for the wet lab comparison (Figure 1) were collected  
150 at the Bermuda Atlantic Time-series Study (31° 40'N 64° 10'W) on expedition BATS 348 on  
151 June 16<sup>th</sup>, 2018, between 01:00 and 05:00 am local time. *In situ* (underwater) large volume  
152 filtration was conducted using submersible pumps to produce replicate biomass samples at a  
153 single depth in the water column for intercomparisons. All filter subsamples are matched for  
154 location, time, and depth. To collect the samples, two horizontal McLane pumps were clamped  
155 together (Figure 1c) and attached at the same depth (80 m) with two filter heads (Mini-MULVS  
156 design) on each pump and a flow meter downstream of each filter head. [This depth was chosen](#)  
157 [to correspond to a depth with abundant chlorophyll and photosynthetic organisms.](#) Each filter  
158 head contained a 142 mm diameter 0.2 µm pore-size Supor (Pall Inc.) filter with an upstream  
159 142 mm diameter 3.0 µm pore-size Supor (Figure 1b, d). Only the 0.2 – 3.0 µm size fraction  
160 was used in this study. The pumps were set to run for 240 min at 3 L per min. Volume filtered  
161 was measured by three gauges on each pump, one downstream of each pump head, and one  
162 on the total outflow (Table S2). Individual pump head gauges summed to the total gauge for  
163 pump 1 (within 1 L; 447 L and 446.2 L), but deviated by 89 L on pump 2 (478 L and 388.9 L).  
164 Given that the total gauge is further downstream, we report the pump head gauges as being  
165 more accurate.

166 The pump heads were removed from the McLane pumps immediately upon retrieval,  
167 decanted of excess seawater by vacuum, placed in coolers with ice packs, and brought into a  
168 fabricated clean room environment aboard the ship. The 0.2  $\mu\text{m}$  pore-size filters were cut in  
169 eight equivalent pieces and frozen at  $-80^{\circ}\text{C}$  in 2 mL cryovials, creating 16 samples per pump  
170 that were co-collected temporally and in very close proximity ( $<1$  m) to each other for a total of  
171 32 samples used in this study (Figure 1d). The 3.0  $\mu\text{m}$  pore-size filters are not included in this  
172 study but are archived for future efforts. The sample naming scheme associated with the  
173 different pumps and pump heads is described in Table S2. Note that pump 1A and 1B samples  
174 accidentally had two 3.0  $\mu\text{m}$  filters superimposed above the 0.2  $\mu\text{m}$  filter, and 1B had a small  
175 puncture in it, although neither of these seemed to affect the biomass collected, presumably the  
176 puncture occurred after sampling was completed.

177 Samples for the [bioinformatic](#) component were collected by the autonomous underwater  
178 vehicle *Clio*. The vehicle and its sampling characteristics were used as previously described  
179 (Breier et al., 2020; Cohen et al., 2023). Specifically, samples Ocean-8 and Ocean-11 were  
180 also collected from the BATS station on R/V *Atlantic Explorer* expedition identifier AE1913 (also  
181 described as BATS validation track BV55 32.75834° N 65.7374° W). The samples were  
182 collected by autonomous underwater vehicle (AUV) *Clio* on June 19th 2019, dive Clio020, with  
183 samples collected at 20 m (Ocean-11) and 120 m (Ocean-8) with 66.6 L and 92.6 L filtered,  
184 respectively, used for this study. [These depths were chosen to reflect the near surface \(high-](#)  
185 [light\) and deep chlorophyll maximum \(low-light\) communities present in the stratified summer](#)  
186 [conditions.](#) These samples were analyzed by 1D DDA analysis using extraction and mass  
187 spectrometry for laboratory 438 [within their laboratory](#) (Tables S5-S7). Sample metadata for  
188 both arms of this intercomparison study and corresponding repository information is provided in  
189 Table S3 and repository links are in the Data Availability Statement.

## 190 2.2 Metagenomic Extraction, Sequencing, and Assembly

191 A metagenomic (reference sequence) database was created for peptide to spectrum  
192 matching (PSMs) for the metaproteomic studies using a 1/8<sup>th</sup> sample split from the exact  
193 sample used in the intercomparison as described above. Samples were shipped on dry ice to  
194 the Naval Research Laboratory in Washington D.C. (USA), where DNA was extracted and  
195 sequenced. Preserved filters were cut into smaller pieces using a sterile blade and placed into a  
196 PowerBead tube with a mixture of zirconium beads and lysis buffer (CD1) from the Dneasy  
197 PowerSoil Pro kit (Qiagen, Hilden Germany). The bead tube with filter sample was heated at  
198 65°C for 10 min then placed on a vortex adapter and vortexed at maximum speed for 10 min.  
199 After sample homogenization/lysis, the bead tube was centrifuged at 16 k x g for 2 min. The  
200 supernatant was transferred to a DNA LoBind tube and processed using the manufacturer's  
201 recommendations. The purified DNA was further concentrated by adding 10 µL 3 M NaCl and  
202 100 µL cold 100% ethanol. The sample was incubated at -30°C for 1 hour, followed by  
203 centrifugation at 16 k x g for 10 min. The supernatant was removed and precipitated DNA was  
204 air-dried and resuspended in 10 mM Tris. DNA concentration was quantified with the Qubit  
205 dsDNA High Sensitivity assay (Thermo Fisher Scientific, Waltham, MA, USA) and DNA quality  
206 was assessed using the NanoDrop (ThermoFisher) and gel electrophoresis. Processing controls  
207 included reagent only and blank filter samples.

208 Sequencing libraries were created from purified sample DNA using the IonExpress Plus  
209 gDNA Fragment Library Preparation kit (Thermo Fisher) for a 200 bp library insert size. No  
210 amplification of the library was required as determined by qPCR using the Ion Library TaqMan  
211 Quantitation Kit. A starting library concentration of 100 pM was used in template generation and  
212 chip loading with the Ion 540 Kit on the Ion Chef instrument prior to single-end sequencing on  
213 the S5 benchtop sequencer.

214 Sequencing used a mix of Ion Torrent and Oxford Nanopore sequencing and resulting  
215 sequencing reads were assembled using SPAdes v. 3.13.1 with Python v. 3.6.8. Following



216 metagenome assembly, contigs smaller than 500 bases were discarded. Open reading frame  
217 (ORF) calling was performed on contigs 500 bps or longer using Prodigal v. 2.6.3 (Hyatt et al.,  
218 2010) run with metagenomic settings as well as MetaGeneMark by submitting to the  
219 MetaGeneMark server ([http://exon.gatech.edu/meta\\_gmhmp.cgi](http://exon.gatech.edu/meta_gmhmp.cgi)) using GeneMark.hmm  
220 prokaryotic program v. 3.25 on August 11, 2019. ORFs called from both programs were  
221 combined and made non-redundant using in-house Python scripts that utilize BioPython v. 1.73.  
222 Non-redundant ORFs were annotated using the sequence alignment program DIAMOND (v 0.9.29)  
223 with the NCBI nr database (downloaded 12/17/2019). ORFs were also annotated with InterProScan  
224 (v 5.29) and with GhostKOALA (Kanehisa et al., 2016) (submitted to server 1/2/2020). Taxonomy  
225 lineages were generated by using the best DIAMOND (Buchfink et al., 2015) hit and pulling lineage  
226 information from NCBI Taxonomy database using BioPython v. 1.73.

### 227 2.3 Proteomic methodologies: Extraction, instrumentation, and *bioinformatics*

228 Some basic protocol stipulations were provided to study participants regarding analytical  
229 conditions to set a uniformity of experimental design. While users were encouraged to use the  
230 extraction method of their preference, constraints on chromatography and mass spectrometry  
231 conditions were set, limiting the number of chromatographic dimensions to one (1D), the total  
232 length of the chromatographic run, the amount of protein injected (as proteolytic digests), and a  
233 single mass spectrometry injection rather than gas phase fraction approaches (Table S4). Each  
234 laboratory group's specific approach is summarized in the supplemental methods, with  
235 extraction in Table S5, and chromatography and mass spectrometry equipment and parameters  
236 in Tables S6 and S7. While there are more sophisticated methods such as two-dimensional  
237 (2D) chromatography and gas phase fractionations that have been demonstrated to provide  
238 deeper metaproteomes (McIlvin and Saito, 2021), these often require specialized equipment  
239 and/or additional instrument time. As a result, the study constraints were provided to ensure a

240 single simple method that all labs could utilize. Laboratories were invited to submit additional  
241 data from more complex analytical setups if they first completed the 1D analyses.

242 ~~Methods used for the informatics intercomparison study are also presented within the~~  
243 ~~Supplemental Materials.~~

244

#### 245 *2.4 Compilation, analysis, and re-analysis of laboratory data submissions*

246 Results from individual laboratories' data submissions were analyzed in two ways as  
247 shown in the flowchart of Figure 1a. First, submitted processed data reports (i.e. PSMs,  
248 taxonomic, functional annotations) were compiled and interpreted. Second, raw data files (i.e.  
249 spectra directly from instruments) from each group were put through a single [bioinformatic](#)  
250 pipeline using SEQUEST HT/Percolator within Proteome Discoverer (Version 2.2.0.388,  
251 Thermo Scientific) and Scaffold (Version 5.2.1, Proteome Software) to isolate variability  
252 associated with bioinformatic processing. [Note that Scaffold ignores the Percolator output from](#)  
253 [Proteome Discoverer when re-running in Scaffold.](#) This re-analysis (*single pipeline re-analysis*  
254 hereon) allowed detailed cross-comparisons of laboratory practices to assess the influence of  
255 the extraction and mass spectrometry components. Specific parameters of the latter included:  
256 parent ~~and fragment~~ [of tolerances of 10ppm were used on all instruments \(all Orbitraps\) for](#)  
257 [fragments tolerances of 0.02 Da or 0.6 Da were used for the instruments with Orbitrap ms2](#)  
258 [instruments and , for ion trap ms2 0.6 Da for ion trap ms2 instruments, respectively, and 0.02](#)  
259 [Da, respectively, with f](#)Fixed and variable modifications of +57 on C (fixed), and +16 on M and  
260 +42 on Peptide N-Terminal (variable) [were used. 0.02 for the instruments with Orbitrap ms2, for](#)  
261 [ion trap ms2 0.6 Da.](#) Peptide and protein FDRs ([false discovery rates](#)) were set to lower than  
262 1.0% using a decoy database, with 1 minimum peptide per protein, and the resulting peptide  
263 FDR was 0.1%. The database used for PSMs was  
264 Intercal\_ORFs\_prodigal\_metagenemark.fasta based on the metagenomic sequencing  
265 described above with 197,824 protein entries. [The protein in this](#) ~~The~~ re-analysis was conducted

266 within Scaffold using total spectral counts and allowing single peptides to be attributed to  
267 proteins. In addition to the total number of protein identifications, the number of protein groups  
268 identified by Scaffold was also provided. Each protein group represented proteins identified with  
269 identical peptides, collapsed into a single protein entry with the highest probability and number  
270 of spectral counts.

271

## 272 *2.5 Data analysis methods*

273 Several analyses were conducted using data from the single pipeline re-analysis. First,  
274 pairwise comparisons of protein identifications were conducted using spectral abundance  
275 reports produced in Scaffold, and loaded, analyzed and visualized in MATLAB (MathWorks Inc).  
276 Two-way (independent) linear regressions were conducted using the script `linfit.m`.  $R^2$  on the  
277 seven datasets were averaged and their standard deviation calculated for shared proteins in  
278 each dataset. Second, a Sørensen similarity (Sørensen, 1948) was calculated where a matrix  
279 was generated that consisted of the unique proteins or peptides identified across all technical  
280 replicates from the various labs with the relative abundance per replicate (% contribution of  
281 each protein/peptide per technical replicate total). The Bray-Curtis dissimilarity pairwise distance  
282 was calculated on this matrix using Python and the SciPy library (v. 1.4.1, (Virtanen et al.,  
283 2020)) and then  $1 - \text{Bray-Curtis dissimilarity}$  was calculated across the matrix to generate the  
284 Sørensen pairwise similarity across all replicates. The resulting similarities per replicate were  
285 clustered and visualized using the `clustermap` function in the Seaborn library (v. 0.10.0,  
286 (Waskom, 2021)). Third, shared peptides and proteins were visualized using Upset plots, using  
287 the R package UpSetR (Conway et al., 2017) to determine the number of unique peptide  
288 sequences and annotated proteins in intersecting sets between all labs, all permutations of lab  
289 subsets, and all lab pairs.

## 290 *2.6. Bioinformatics Intercomparison Methods*

Formatted: Font: Italic

Formatted: Font: Italic

291 ~~The methods used for the bioinformatics intercomparison study are also presented within the~~  
292 ~~Supplemental Materials~~ are described by each laboratory using their unique three-digit identifier  
293 code. All laboratories used the metagenomic database generated in the laboratory study (see  
294 Section 2.2).

295 **Lab 109:** The raw files were searched against the metagenomic database employing a 2 round  
296 search using PEAKS Studio X. The initial database search was performed to focus the  
297 metagenomic database for protein sequences with peptide sequence matches at 5% FDR. The  
298 focused database was further used for a second round search, which allowed a parent mass  
299 error tolerance of 10.0 ppm and a fragment mass error tolerance of 0.6 Da. The search  
300 considered up to 3 missed cleavages, carbamidomethylation as fixed and methionine oxidation  
301 and N-terminal acetylation as variable modifications. The cRAP protein sequences  
302 (<http://ftp.thegpm.org/fasta/cRAP/>) were included as contaminant database. Finally, PSMs were  
303 filtered for 1% FDR and annotated with taxonomic lineages (obtained from the metagenomic  
304 experiments). Non-unique peptide matches were annotated with the LCA of the respective  
305 lineages.

306 **Lab 321:** SearchGUI (Galaxy Version 3.3.10.1) was used to search using multiple search  
307 algorithms (X!Tandem, MS-GF+ and Comet). For each search algorithm, Precursor Tolerance  
308 of 10.0 ppm, Fragment Ion Tolerance of 0.6 Da and trypsin was used as an enzyme for  
309 proteolytic cleavage. Searches were performed allowing for two missed cleavages fixed  
310 modification of Carbamidomethylation at cysteine and Variable Modifications of Acetylation of  
311 protein N-term and Oxidation of Methionine. PeptideShaker (Version: 1.16.36) was used to filter  
312 peptides with the length of 8-50 aas and a precursor m/z tolerance of 10.0 ppm. Detected  
313 peptide-spectral matches, peptides and proteins were reported at 1% global FDR. All of the  
314 analysis was performed within Galaxy platform.

315 **Lab 321:** MaxQuant (Galaxy version 1.6.17.0+galaxy3) was used to search the datasets. A  
316 fixed modification of carbamidomethylation at cysteine and variable mmodifications of

Formatted: Indent: First line: 0"

317 acetylation of protein N-term and oxidation of methionine was applied along with allowing for  
318 two missed cleavages. The detection peptides and proteins were reported at 1% FDR.  
319 Lab 362: The raw files were converted using ThermoRawFileParserGUI (version 1.4.1) to peak  
320 lists (.mgf files) using "native Thermo library peak picking" as the peak picking option and  
321 "Ignore missing instrument properties" as the error option. The peak lists (.mgf files) obtained  
322 from MS/MS spectra were identified using X! Tandem version X! Tandem (Vengeance version  
323 2015.12.1) using SearchGUI version 4.1.0. Here, the parameters provided and suggested by  
324 the study were used: tolerances of 10 ppm for MS1 and 0.6 Dalton for MS/MS; dynamic  
325 modifications: oxidation of M, and acetyl on N-terminus; static modifications: carbamidomethyl  
326 of C. Identification was conducted against a concatenated target/decoy database of the  
327 provided database.  
328 The X!Tandem files were used as input in MS<sup>2</sup>ReScore  
329 (<https://github.com/compomics/ms2rescore>), a machine learning-based post-processing tool  
330 that improves upon Percolator rescoring of peptide-to-spectrum matches (PSMs). Here, the  
331 search engine-dependent features of Percolator were appended with MS<sup>2</sup> peak intensity  
332 features by comparing the PSM with the corresponding MS<sup>2</sup>PIP-predicted spectrum. All  
333 reported MS<sup>2</sup>ReScore PSM identifications have a q-value < 0.01. No protein grouping algorithm  
334 was applied, and all identified taxa and functions are extracted from the provided database.  
335 Lab 458: The Proteome Discoverer 2.5 platform was used (SequestHT + Percolator (MPS)).  
336 Fully tryptic peptides with a minimum length of 6 peptides and a maximum of 2 missed  
337 cleavages were required. Precursor Tolerance of 10.0 ppm, Fragment Ion Tolerance of 0.6 Da,  
338 carbamidomethylation as fixed and methionine oxidation was set as a variable modification. Filtering  
339 was performed at a 1% PSM- and peptide-level FDR. The MaxQuant contaminant list was used as  
340 a contaminant database.  
341 Lab 501: We first appended the database with a set of common contaminants (Global  
342 Proteome Machine Organization common Repository of Adventitious Proteins). Then, we used

Formatted: Font: Bold

343 MSGF+ (Kim and Pevzner, 2014) to match mass spectra with peptide sequences, with cysteine  
344 carbamidomethylation as a fixed modification, and methionine oxidation, glutamine modified to  
345 pyro-glutamic acid, deamidated asparagine, and deamidated glutamine, as variable  
346 modifications. Peptides were searched for with a Target-Decoy approach, with a 1% false  
347 discovery rate at the peptide spectrum match level. For spectral counts, we summed MS2  
348 spectra that identified a peptide, and normalized all spectral counts to the total spectral counts  
349 per sample. Proteins were quantified using the median spectral count for all proteotypic  
350 peptides (those peptides which uniquely correspond to a protein), specifically using the  
351 OpenMS tool ProteinQuantifier. This approach requires at least one proteotypic peptide, but if  
352 more are identified, those peptides are also used for quantification.

353 Lab 828: The raw files were analyzed using Thermal proteome discover. MS/MS spectrums  
354 were searched against provided database using SEQUEST-HT engine. MS/MS spectra  
355 searches were performed as follows: precursor ion tolerance of 10.0 ppm; fragment ion  
356 tolerance of 0.6 Da; carbamidomethyl cysteine was specified as fixed modification, whereas  
357 oxidation (M), deamidation (N/Q), and N-terminal protein acetylation were set as variable  
358 modifications. Trypsin was specified as the proteolytic enzyme, allowing for two missed  
359 cleavages. Percolator-based scoring was chosen to improve the discrimination between correct  
360 and incorrect spectrum identifications, learning from the results of a decoy and target database;  
361 settings were as follows: maximum delta Cn, 0.05; strict false-discovery rate of 0.01 and  
362 validation based on q values.

363 Lab 902: SEQUEST-HT was used within Proteome Discoverer 2.2 using the following settings:  
364 maximum missed cleavage 2, minimum peptide length 6, maximum peptide length 122,  
365 precursor mass tolerance 10ppm, fragment mass tolerance 0.6 Dalton; dynamic modifications:  
366 M oxidation, acetyl on N-terminus; static modifications: C carbamidomethyl. Percolator PSM  
367 validator (within Proteome Discoverer) with following settings: maximum Delta Cn 0.05, target  
368 FDR strict 0.01, target FDR relaxed 0.05, validation based on PEP. Scaffold 5.0 used to analyze

369 [Proteome Discoverer generated files with following settings: scoring system: prefiltered mode;](#)  
370 [protein grouping: standard experiment wide protein grouping; protein threshold 1.0% FDR;](#)  
371 [peptide threshold 0.1% FDR; minimum number of peptides 1.](#)  
372 [Lab 932: Mass spectrometry data were transformed from Thermo RAW format \(version 66\) to](#)  
373 [mzML and Mascot Generic \(MGF\) formats using ThermoRawFileParser \(version 1.2.0,](#)  
374 [Hulstaert et al., 2020\). Experimental metadata were extracted from mass spectrometry data](#)  
375 [using the MARMoSET program \(Kiweler et al. 2019\). Mascot Server \(version 2.6.2, Matrix](#)  
376 [Science, LTD\) software performed peptide-spectrum matching between experimental data and](#)  
377 [a reference sequence database. Reference sequences included a total of 197,824 predicted](#)  
378 [protein-coding ORFs from a metagenome assembly. Peptides matching an in-house curated](#)  
379 [inventory of contaminant protein sequences, mass standards, and proteolytic enzyme](#)  
380 [sequences were removed from the results. Mascot search parameters included the following](#)  
381 [settings: +10.0 ppm monoisotopic precursor mass tolerance; +0.6 Da monoisotopic fragment](#)  
382 [ion tolerance: one fixed modification \(+57 to C residues\); two variable modifications \(+16 to M](#)  
383 [residues, +42 to peptide amino-termini\); digestion enzyme trypsin; two missed cleavages;](#)  
384 [peptide charges +2-+7; and instrument type: electrospray ionization coupled to fourier-transform](#)  
385 [ion cyclotron resonance \(ESI-FTICR\). Mascot search results containing peptide-spectrum](#)  
386 [matches \(PSMs\) were exported for downstream data analysis. Scaffold Q+S \(version 4.8.9\) was](#)  
387 [used to validate MS/MS-based peptide- and protein-level peptide-spectrum matches \(PSM\) with](#)  
388 [the Peptide Prophet algorithm. Mascot PSM data were imported into Scaffold Q+S with the](#)  
389 [following settings specified: quantitative metric: spectrum counting; scoring system: use legacy](#)  
390 [Peptide Prophet scoring \(high mass accuracy\); protein grouping: use standard experiment-wide](#)  
391 [grouping; optional loading steps: pre-compute false discovery rate \(FDR\) thresholds; and, use](#)  
392 [local gene ontology \(GO\) annotations \(UniProt GO annotation data retrieved 25 JUN 2020\).](#)  
393 [Scaffold Q+S identification criteria were set at greater/equals >99.9% probability by the Peptide](#)

394 [Prophet algorithm \(Keller et al. Anal. Chem. 2002.\)](#) and >99.9% probability by the Protein  
395 [Prophet algorithm \(Nesvizhskii et al., Anal. Chem. 2003\)](#) with >2 peptides at the protein level.  
396 **Lab 957:** MSFragger 3.3 searches were performed with FragPipe 16.0 and Philosopher 4.0.0. A  
397 concatenated target/reverse database was searched with a 50 PPM precursor and 0.4 Da  
398 fragment mass tolerance. Automatic mass calibration and parameter optimization was enabled  
399 and precursor mass errors for up to +2 neutrons were considered. Peptide candidates were  
400 generated from database protein sequences assuming tryptic digestion, allowing for up to one  
401 missed cleavage. Peptides were required to have between 8-50 amino acids and range from  
402 500 to 5000 m/z. Cysteines were assumed to be fully carbamidomethylated, and peptides were  
403 searched considering variable n-terminal pyroglutamic acid formation and methionine oxidation.  
404 PeptideProphet was used for FDR validation with the following default options: "--decoy probs",  
405 "--ppm", "--accmass", "--nonparam", and "--expectscore", which allow for additional high-mass  
406 accuracy analysis and non-parametric distribution fitting. ProteinProphet was used for protein-  
407 level FDR validation with the following default option: "--maxppmdiff 2000000". Filtering was  
408 performed using a 1% peptide-level and a 1% protein-level FDR threshold.

Formatted: Indent: First line: 0"

### 410 3. Results

#### 411 3.1 Experimental Design

412 This ocean metaproteomic intercomparison consisted of two major [components](#) activities:  
413 a laboratory component, where independent labs processed identical ocean samples  
414 simultaneously collected from the North Atlantic Ocean (Fig. 1a, see Section 2.1), and a  
415 subsequent [bioinformatic](#) component. Participating institutions and persons at those institutions  
416 are listed in Table S1, with all participants also listed as co-authors. Both arms of the study were  
417 conducted under blinded conditions, where correspondence with participants was conducted by  
418 an individual not involved in either study, and submitted results and data were anonymized prior



419 to sharing with the consortium. Within both arms of the study, participants were provided the  
420 location of the study site and metadata about the sampling locations, time and depth at the  
421 onset of the study. The laboratory study involved two biomass-laden filter slices collected from  
422 the North Atlantic Ocean Bermuda Atlantic Time series Study site at 80m depth being sent to  
423 each participating group for protein extraction, mass spectrometry, and bioinformatic analyses  
424 (see Section 2.1 below). This depth was chosen to correspond to a depth with abundant  
425 chlorophyll and associated photosynthetic organisms. The bioinformatic effort was independent  
426 of the laboratory effort and involved the distribution and bioinformatic analysis of two  
427 metaproteomic raw data files generated from samples also from the North Atlantic Ocean upper  
428 water column BATS station (20m and 120m depths, see Section 2.1). These depth were chosen  
429 to reflect the near surface (high-light) and deep chlorophyll maximum (low-light) communities  
430 present in the stratified summer conditions. These files were distributed after labs had submitted  
431 their laboratory extracted raw data files. The raw files from the bioinformatic study were distinct  
432 from the samples used in the laboratory intercomparison study to avoid any biases from groups  
433 that analyzed those samples previously. Submitted results from both components were  
434 anonymized and assigned three-digit lab identifiers generated randomly with laboratory and  
435 bioinformatic results from the same lab being assigned distinct identifiers.

Formatted: Font: Not Italic

Formatted: Font: Not Italic

436 We report results for two study components: Part 1 (Section 3.2) involves the data  
437 generation intercomparison of distributed subsamples from the North Atlantic Ocean (Fig. 1;  
438 Section 2.1). Part 2 (Section 3.3) was an bioinformatic intercomparison, where metaproteomic  
439 raw files were shared with participants and processed results were submitted. Both components  
440 were conducted as blinded studies, where each dataset was assigned a three digit randomly  
441 generated identifier, with those identifiers used throughout the Results and Discussion.

442

443 *3.2 Mass Spectrometry Data Generation Intercomparison*

444           Nine laboratories submitted raw and processed datasets from the analysis of the  
445 distributed Atlantic Ocean field samples (Table S1). The processed data submissions were  
446 heterogeneous in output formats, statistical approaches, and parameter definitions. Because of  
447 the challenges of comparing data derived from different types of statistical approaches used for  
448 peptide and protein identification and inference, as well as the varying output formats from  
449 various software packages, the user-generated data submissions were difficult to compile and  
450 compare, resulting in variability in the number of identifications depending on the statistical  
451 approaches and thresholds applied. These results are further discussed in the Supplemental  
452 Section (Figure S1, Table S8). Despite these challenges, an average of 7142 +/- 2074 peptides  
453 were identified across the pairwise comparisons (Figure S1c) representing 20% of the 35,715  
454 total unique peptides detected across all labs. Together these findings imply a  
455 consistency of peptide identifications across participants. The variability in proteome depth  
456 reflected the combination of differing parameters employed by software and laboratory  
457 approaches.

458           To remove this variability associated with user-selected bioinformatic pipelines, a single  
459 pipeline re-analysis of the submitted raw mass spectral data was conducted. Raw data files  
460 were processed together within a single bioinformatic pipeline consisting of SEQUEST-HT,  
461 Percolator, and Scaffold software and evaluated to a false discovery rate threshold of < 0.1% for  
462 peptides and 1.0% for proteins % (see see Section 2.4). Two datasets were found to have had  
463 issues during extraction and analysis that affected the results in both processed and raw data  
464 (Labs 593 and 811; Table S8). Notably these two laboratories differed from the others in that  
465 they did not use SDS as a protein solubilizing detergent (Table S5). This likely resulted in  
466 inefficient extraction of the bacteria that dominated the sample biomass (e.g. picocyanobacteria  
467 and *Pelagibacter*) embedded within the membrane filter slices. Further examination showed  
468 polyethylene glycol contamination of one dataset (Lab 811) and low yield from sample  
469 processing and extraction from the other (Lab 593). As a result, those datasets were not

Formatted: Font: Not Italic

470 included in the single pipeline re-analysis. The standardized pipeline included calculations of  
471 shared peptides and proteins, quantitative comparisons, and consistency of taxonomic and  
472 functional results.

473 The total number of peptide and protein identifications and PSMs in the single  
474 bioinformatic pipeline analysis varied by laboratory (Table S9), with unique peptides ranging by  
475 more than a factor of 3 from 3,354 to 16,500, and with 27,346 total unique peptides identified  
476 across laboratories. This variability was likely due to different extraction, chromatographic, and  
477 mass spectrometry ~~hardware and parameters employed approaches~~ used by each laboratory,  
478 resulting in a varying depth of metaproteomic results. Yet, as with the user-submitted results,  
479 there was considerable overlap in identifications between all datasets. An intersection analysis  
480 found the numerous shared peptides between all combinations of laboratories, with 1,395  
481 peptides shared between all seven laboratory datasets (Figure 2a). Laboratories with deeper  
482 proteomes shared numerous peptides, for example the two laboratories with the most  
483 discovered unique peptides shared ~3000 peptides between them, implying that shared  
484 peptides is a useful metric for intercomparability. They also had the largest numbers of peptides  
485 that were not found by any other labs (3617 and 2819, respectively). The fourth largest  
486 intersection size (1395) represented the unique peptides discovered by all labs. Beyond that  
487 there were 12 different groupings of peptides that were shared among at least four laboratories.  
488 Consistent with this, 3-way Venn diagrams of labs 135, 209 and 438 had an intersection of 2398  
489 peptides, labs 652, 729, and 774 ~~shared~~ 3016 peptides, and labs 127, 135, and 309  
490 shared 2304 peptides (Figure 2d).

491 A similar analysis was conducted at the protein level, where the number of proteins  
492 identified ~~for each sample based on peptide mapping to the metagenome database~~ (see  
493 Section 2. Methods), ~~identified~~ 8,043 ~~total~~ unique proteins in total across all ~~seven~~ laboratories,  
494 ~~with and~~ 1,056 proteins of those observed in shared amongst those laboratories all seven labs  
495 (see as shown in the 7-way Venn diagram in (Figure 2c). Three-way Venn diagram comparisons

496 among labs 135, 209 and 438 had an intersection of 1,254 proteins, and labs 652, 729, and 774  
497 shared 1,925 proteins (data not shown).

498 Optional deeper metaproteome results were submitted by three laboratories using either  
499 a long gradient of 12 hours or 2 dimensional chromatographic methods (Table S10). The  
500 number of discovered peptide and protein identifications were higher in each case, with as  
501 many as 18477 unique peptides and 7765 protein identifications from an online 2-dimensional  
502 chromatographic analysis from a 5 µg single injection.

503 The mapping of identified peptides to protein sequences forms the basis for protein  
504 identifications in the form of DDA bottom-up proteomics employed here. The relationship  
505 between peptides and protein identification was explored in Figure 3 and found to be correlated  
506 by two-way linear regression with  $R^2$  values of 0.97 and 0.98 for total protein identifications and  
507 protein groups, respectively. Together, the fact that there is a linear relationship between  
508 peptides and proteins across all laboratories (including labs employing deeper methods) could  
509 imply that the number of protein identifications has not begun to plateau and reached  
510 'saturation', likely due to the immense biological diversity and abundance of lower abundance  
511 peptides within these samples. This approach has some similarities to rarefaction curves used  
512 in metagenomic sequencing to determine if the majority of species diversity has been sampled,  
513 although in this case number of peptides used as a metric for sampling depth instead of  
514 additional number of DNA sequencing samples typically used for rarefaction curves. This  
515 indicated that with deeper depth of analysis by some laboratories, there was no fall off in the  
516 increase in protein identifications that might be attributed to additional peptides mapping to  
517 already discovered protein sequences. In addition, the 2D and long gradient additional analyses  
518 conducted by several laboratories fell upon this line consistent with this “more peptides – more  
519 proteins” observation, implying more room for improvements in depth of metaproteomic  
520 analyses.

521 A quantitative analysis of spectral counts from the wet lab re-analysis showed broad  
522 coherence among the seven laboratories. Pairwise comparisons of protein spectral counts were  
523 conducted for each of the seven labs against the other six (visualized in a 7x7 matrix, with  
524 duplicate comparisons removed (e.g., A vs B and B vs A)), where each data point reflects the  
525 spectral counts for a protein shared between laboratories (Figure 4a). When a dataset was  
526 compared with itself a unity line of datapoints was observed along the diagonal axis as  
527 expected. Two-way linear regressions were conducted on each of these pairwise comparisons.  
528 The slopes ranged from 0.33 to 5.5 (Figure S2), implying a varying dynamic range in spectral  
529 counts across laboratories, likely due to variations in instrument parameterizations selected by  
530 each laboratory, and consistent with the lack of normalization between laboratories. The  
531 coefficient of determination  $R^2$  values from 0.43 to 0.8473 with an average of 0.63 +/- 0.11,  
532 showing inged coherence among results for these large metaproteomic datasets (Figure 4b, Table  
533 S12). To provide a sense of coherence of each laboratory to the others, the  $R^2$  values of a lab  
534 against the other six laboratories were averaged and the standard deviation calculated. All of  
535 these average  $R^2$  values were higher than 0.5, which showed overall quantitative consistency  
536 despite the size and complexity of these datasets (Figure 4d).

537 A comparative taxonomic and functional analysis was also conducted using a single  
538 bioinformatic pipeline (see metagenomic sequencing methods for annotation pipeline). Lowest  
539 common ancestor (LCA) analysis of peptides identified from datasets from seven laboratories  
540 showed consistent patterns of taxonomic distribution using the MetaTryp package (Figure 5a;  
541 (Saunders et al., 2020). Cyanobacteria and alphaproteobacteria were the top two taxonomic  
542 groups in all laboratory submissions, consistent with the abundant picocyanobacteria  
543 *Prochlorococcus* and the heterotrophic bacterium *Pelagibacter ubique* known to be dominant  
544 components of the Sargasso Sea ecosystem (Sowell et al., 2009; Malmstrom et al., 2010). For  
545 example, *Prochlorococcus* is consistently present between  $10^4$  and  $10^5$  cells per milliliter -in this  
546 region and has been observed to contribute to carbon export from the euphotic zone (Casey et

547 al., 2007). *Pelagibacter* cells can also be in excess of 10<sup>5</sup> cells per milliliter at the BATS North  
548 Atlantic location (Carlson et al., 2009). These results are broadly similar to the representation of  
549 phyla within the metagenome annotations, where Proteobacteria (including *Pelagibacter*) and  
550 Cyanobacteria (including *Prochlorococcus* and *Synechococcus*) were major components,  
551 although Bacteroidetes (including Flavobacteria) are more prevalent in the metagenome  
552 annotations than in the metaproteome. Some differences may also be due to the incorporation  
553 of protein abundances in Fig 5a, versus simple taxonomic attribution of non-redundant  
554 assembled open reading frames in the metagenome analysis, as well as the use of multiple  
555 sequencing platforms and gene calling algorithms (Section 2.2, Figure S4).

556 Similarly, KEGG Orthology group (KO) analysis of those datasets also showed highly  
557 similar patterns of protein functional distributions across laboratories (Figure 5b). Notably the  
558 PstS phosphate transporter protein from *Prochlorococcus* was the most abundant functional  
559 protein in all datasets, consistent with observations of phosphorus stress in the North Atlantic  
560 oligotrophic gyre and its biosynthesis in marine cyanobacteria (Scanlan et al., 1997; Coleman  
561 and Chisholm, 2010; Ustick et al., 2021). These findings demonstrate the reproducibility in the  
562 primary functional and taxonomic conclusions from the metaproteome datasets. Finally, a  
563 Sørensen similarity analysis of the 1,000 proteins with highest spectral counts revealed 70–80%  
564 similarities between most laboratory groups in the data re-analysis (Figure 6). When conducted  
565 on the Similarity analyses on the full dataset (with all peptides and proteins), the Sørensen  
566 similarity analyses showed revealed lower similarity at the peptides had lower similarity than  
567 level than the proteins level, implying variability in peptide identification is ameliorated when as it  
568 is aggregated to the protein level (Figure S3).

### 570 3.3. *Bioinformatics Data Analysis Intercomparison*

571 Two metaproteomic raw files were provided to intercomparison participants and were  
572 searched with each laboratory's preferred database searching bioinformatics pipeline. The

Formatted: Font: Italic

573 samples that generated the data for these files were collected by autonomous AUV *Clio* during  
574 a single dive at the Bermuda Atlantic Time-series Study Station (Breier et al., 2020), and were  
575 distinct from the samples associated with the laboratory intercomparison component. However,  
576 they were also from the North Atlantic Ocean, allowing the same metagenomic database to be  
577 used. This database was not collected simultaneously with the bioinformatics samples, so it was  
578 not as representative as that used in the laboratory intercomparison. However, the BATS study  
579 region is known to maintain similar major taxonomic composition throughout the year (e.g.,  
580 *Prochlorococcus* and SAR11, see discussion in Section 3.2), hence enabling many protein  
581 identifications. This bioinformatic study component was not launched until after the laboratory-  
582 based intercomparison submission deadline to avoid influencing that part of the study by  
583 sharing similar raw data. Samples were named Ocean 8 and Ocean 11 and were taken from  
584 120 m and 20 m depths, respectively.

585 The bioinformatic intercomparison involved 10 laboratories utilizing 8 different software  
586 pipelines including the PSM search engines: SEQUEST, X!Tandem, MaxQuant, MSGF+,  
587 Mascot, MSFragger, and PEAKS (Table S11, see Methods Section 2.6). As with the user  
588 supplied laboratory results, the results were challenging to compile due to different types of data  
589 outputs, approaches used in protein inference, and statistical approaches applied within each  
590 pipeline. Unique peptide discoveries served as a useful base unit of comparison that were less  
591 subject to these comparison challenges. The number of peptides ranged from 1724 to 6369 in  
592 Ocean 8 and 3019 to 8288 in Ocean 11 (Figure 7; Table S11). The differences in the number of  
593 peptides was likely due to parameters used in software, for example, laboratory 932 had the  
594 lowest number of peptides identified in both samples, but also used a highly stringent 99.9%  
595 probability cutoff that likely influenced this result.

596

#### 597 **4. Discussion**

Formatted: Font: Italic

598 4.1 Assessment of Ocean Metaproteomics Reproducibility

599 Given the ~~relatively~~ recent establishment of ~~complex ocean~~ metaproteomic techniques  
600 ~~as well as their methodological complexity~~, intercomparisons ~~of methods~~ are ~~valuable~~ ~~important~~  
601 in demonstrating ~~their the~~ suitability ~~of metaproteomic analyses in for~~ ocean ecological and  
602 biogeochemistry studies. Synthesizing the results of the laboratory and mass spectrometry  
603 blinded intercomparison study (Section 3.2) processed with a single ~~bio~~informatic pipeline  
604 (Section 2.4), we observed consistent reproducibility with regards to three attributes of ocean  
605 metaproteomics analyses: 1) the identity of discovered peptides and proteins (Fig. 2), 2) their  
606 relative quantita~~ive abundances~~~~tion~~ (Figs. 4 and 6), and 3) the taxonomic and functional  
607 assignments within intercompared samples (Fig 5). With over 1000 proteins identified across  
608 seven laboratories and Sørensen similarity indexes typically higher than 70–80% (Fig. 6), the  
609 results ~~unambiguously~~ demonstrate consistent detection and quantitation of major proteins in  
610 the sample. ~~Together~~ ~~These~~ results provide confidence that multiple laboratories can generate  
611 reproducible results describing the major proteome composition of ocean microbiome samples  
612 ~~to, and in doing so can~~ assess their functional ~~composition~~ and biogeochemical ~~activity~~  
613 ~~significance of these complex microbial communities.~~

614 While there is good agreement, this congregation of data allows further exploration of  
615 the influence of methods on the results. In particular, as mentioned above the range of pairwise  
616 comparisons had correlation coefficients ranging from 0.43 to 0.84, with most values falling  
617 between 0.6 and 0.8 (Figure 4b and 4e; Table S12). This average of all correlation coefficients  
618 described above (0.63 +/- 0.11) implied good reproducibility between laboratories in general.  
619 We can explore what might have influenced the variability and lower range of coefficients. The  
620 correlation coefficients of lab 209 had two of the three R<sup>2</sup> values below 0.499 in pairwise  
621 comparisons (0.431 and 0.475), yet also had values that ranged from 0.61 to 0.70. Why would  
622 this variability exist? Lab 209 's methods differed from other labs in several ways: they used the



623 oldest and slowest instrument of the group (Thermo Orbitrap Elite), used CID instead of HCD for  
624 fragmentation and rapid scan mode, and used an unusually long column of 200cm to  
625 compensate for the older instrument (Table S6). As a result, lab 209 had the lowest number of  
626 peptide (3354) and protein (1586) ID's of the seven labs (Table S9), which was several fold  
627 lower than the lab with the highest number and reduced the number of shared peptides across  
628 all laboratories. In pairwise comparisons, lab 209 had the lowest number of shared peptides at  
629 an average of 1304. Interestingly however, lab 209 did not have the lowest number of total  
630 spectral counts (63198), being close to the average (70843 +/- 27455), implying that more  
631 abundant peptides were detected relative to rarer ones.

632 We initially suspected the lower R<sup>2</sup> values in pairwise comparisons with lab 209 may  
633 have been related to comparisons to laboratories with similarly lesser peptide depth, but this  
634 was not the case: the two lowest correlation coefficients for lab 209 were with laboratories 135  
635 and 774 (the 0.431 and 0.475 values), the latter of which had the highest number of peptide  
636 identifications. The answer for this difference in quantitative values maybe within the selection of  
637 parameters used to sample peptide peaks: Both lab 135 and 774 used 60 second dynamic  
638 exclusion, whereas the other 5 labs used dynamic exclusions between 10 and 30 seconds in  
639 length (Table S7). This higher dynamic exclusion likely contributed to providing greater peptide  
640 discovery depth, but at the cost of quantitative consistency with other laboratories, since this  
641 parameter selects against repeat counting of abundant peaks and would reduce spectral counts  
642 of the more abundant peptides that lab 209 was detecting. This result demonstrates the  
643 influence of the mass spectrometer parameters in quantitative reproducibility when using global  
644 proteomic DDA mode.

645 *4.2 Metrics in metaproteomics: Core versus rare "long tail" proteins*

Formatted: Font color: Text 1

646 While abundant proteins were consistently detected across seven laboratories'  
647 submissions, there was substantial variability in the less abundant proteins (Fig. 2). This is  
648 evident in Figure 8, where most of the 1063 proteins across seven laboratories in the re-  
649 analysis were in the upper half of proteins when ranked by abundance. This simultaneous  
650 consistency in abundant proteins and diversity in rare proteins (and their respective peptide  
651 constituents) was likely a result of several factors, in the study design and execution. First, the  
652 intercomparison experimental design stipulated 1D chromatography in order to provide  
653 straightforward comparisons that all laboratories could accomplish. This contributed to study  
654 consistency, but also resulted in lesser proteome depth compared to more elaborate methods  
655 such as 2D chromatography and gas phase fractionation commonly in use. Second, the sample  
656 complexity of ocean metaproteomes has been shown to be enormous, with a far greater  
657 number of low abundance peptides present than HeLa human cell lines (Saito et al., 2019). The  
658 combined effect of these factors meant that, while laboratories were able to detect abundant  
659 proteins consistently, there was considerable stochasticity associated with the detection of less  
660 abundant peptides resulting in a long tail of discovered lower abundance proteins. This is  
661 evident in Figure 8, where most of the 1063 proteins across seven laboratories in the re-  
662 analysis were in the upper half of proteins when ranked by abundance.

663 Mass spectrometer settings such as dynamic exclusion, chromatography conditions, and  
664 variation in sample preparation methods all likely contributed to this stochastic variability in rare  
665 peptide detection among laboratories. Moreover, while all participating laboratories used  
666 Thermo orbitrap mass spectrometers, there were seven variants of instrument model, including  
667 some with Tribrid multiple detector capability (Table S6). While testing other mass spectrometry  
668 platforms is of interest, this trend of community orbitrap usage in this study is consistent with the  
669 broader proteomics community, where currently 9 of the top 10 instruments used in  
670 ProteomeXchange consortium repository data submissions utilize orbitraps as of the manuscript

671 submission date (Deutsch et al., 2019). When conducting analysis of environmental samples,  
672 choices can be made about instrument setup and parameters based on the scientific objectives,  
673 for example if maximal proteome depth or robust quantitation while using a discovery approach  
674 is desired. Future intercalibration efforts enlisting more sensitive metaproteomic methods such  
675 as 2D-chromatography (McIlvin and Saito, 2021), more sensitive instruments (Stewart et al.,  
676 2023), and other emerging methods can greatly improve detection and quantitation of rarer  
677 proteins in metaproteomes, allowing exploration of the depths of state-of-the-art capabilities  
678 rather than our present emphasis on interlaboratory consistency. Moreover, the development  
679 and adoption of best practices in sample collection, extraction, chromatographic separation,  
680 mass spectrometry analyses, and [bioinformatic](#) approaches will contribute to interlaboratory  
681 consistency.

682         Despite the inter-laboratory variability in the detected sets of rarer peptides and proteins,  
683 we interpret these to be largely robust identifications. The stringent 0.1% peptide-level FDR  
684 threshold we use here is determined by scoring decoys: reverse sequenced peptides that are  
685 not in our samples. Peptide assignments to these decoys model the score distribution of all  
686 incorrect peptide-spectrum matches (PSMs) in our study such that FDRs can be estimated in an  
687 unbiased way for each laboratory. However, these estimates are complicated by subtle  
688 sequence diversity within a population's proteome, which is typically not considered by  
689 proteomics software designed to analyze single species (Schiebenhoefer et al., 2019). This  
690 diversity within metaproteomic samples results in the presence of highly similar peptides with  
691 nearly identical precursor masses that produce many of the same b- and y-ions, and this  
692 similarity is not well modeled by decoy peptides. The influence of microdiversity on  
693 metaproteomics FDR estimation using strain-specific proteogenomic databases is an important  
694 area of future exploration (Wilmes et al., 2008).

695 [4.4.4.3 Bioinformatics Intercomparison Assessment](#)

696 The discovery of peptide constituents of proteins within a complex ocean metaproteomic  
697 matrix was successful across all software packages tested (Figure 7), where the metric for  
698 success is a comparable number of peptide identifications. This is a notable finding due to the  
699 highly complex mass spectra, large number of chimeric peaks present (Saito et al., 2019), and  
700 large database sizes involved in ocean metaproteomes. To our knowledge, some of these  
701 software packages had not yet been applied to ocean metaproteomes. There was also  
702 variability associated with the stringency of statistical parameters employed, which points to the  
703 challenges in assembling datasets from multiple laboratories with different depth of proteome  
704 identifications.

705 Despite the success of this intercomparison component across software packages, there  
706 is likely considerable room for improvement in the future. As mentioned previously, ocean  
707 samples are highly complex and there are likely additional peptides that remain unidentified  
708 using current technology, due to low intensity peaks and co-elution with other peptides resulting  
709 in the chimeric spectra. Significant improvements in depth of analysis can be achieved through  
710 increased chromatographic sample separation and optimized (or alternative) mass spectrometry  
711 data acquisition strategies. Yet there is room for [bioinformatic](#) improvements as well: most DDA  
712 database searching algorithms are unable to identify multiple peptides within a single  
713 fragmentation spectrum. Moreover, when in DDA collection mode mass spectrometry software  
714 typically does not isolate and fragment peptides that cannot be assigned a charge state, which  
715 is a common occurrence for the low abundance peaks within ocean samples. As a result, there  
716 is considerable room for improvements in [bioinformatic](#) pipelines to discover additional peptides.  
717 Although the application of data independent approaches (DIA) to oceanographic  
718 metaproteomics analysis ~~has been is currently~~ limited ([e.g.](#) Morris et al., 2010), the systematic  
719 nature of ion selection and fragmentation allows for a greater number of low abundant peptides  
720 to be quantified. ~~By avoiding the need to select precursor ions for fragmentation, DIA methods~~

721 ~~have the promise to identify some of these rarer peptides, when assuming~~ enough ions can be  
722 isolated to produce robust MS2 spectra, ~~as the wider isolation windows often used in DIA will~~  
723 ~~dilute precursor ions within ion traps.~~

724 4.5.4.4 Lessons Learned and Future Efforts in Ocean Metaproteomic Intercomparisons and  
725 Intercalibrations

726 As the first interlaboratory ocean metaproteomics study, we chose to describe this study  
727 as an intercomparison rather than an intercalibration and it served as a vehicle with which to  
728 assess the extent of reproducibility. There were several lessons learned that can be  
729 summarized here. These include the efficacy of a SDS detergent and heat treatment in lysing  
730 and solubilizing marine microbial cells embedded on membrane filters, the significant problem  
731 of data intercomparability between PSM software outputs and need for data output  
732 standardization, and the influence of different hardware capabilities (Orbitrap generation) and  
733 their parameter settings such as dynamic exclusion on proteome depth and quantitative  
734 comparisons of spectral counts. ~~As mentioned above, The development of best practices~~  
735 associated with sample collection, extraction, and analysis would be valuable, while also  
736 encouraging methodological improvements and backward compatibility through the use of  
737 reference samples.

Formatted: Font: Not Italic

738 Future intercalibration efforts could aim to further assess and improve upon the level of  
739 accuracy, reproducibility, and standardization of ocean metaproteome measurements. ~~As~~  
740 ~~mentioned above, development of best practices associated with sample collection, extraction,~~  
741 ~~and analysis would be valuable, while also encouraging methodological improvements and~~  
742 ~~backward compatibility through the use of reference samples. In particular, Alternative modes~~  
743 of data collection and quantitation could also be tested in future interlaboratory comparisons,  
744 including parallel reaction monitoring mode (PRM), multiple reaction monitoring mode (MRM),

745 quantification using isotopic labeling or tagging, and DIA methods. PRM and MRM methods  
746 allow sensitive targeted measurements of absolute quantities of peptides (e.g. copies per liter of  
747 seawater in the ocean context). As many 'omics methodologies applied in environmental  
748 settings operate in relative abundance modes, adding the ability to measure absolute quantities  
749 would be particularly valuable for comparisons of environments across space and time.  
750 Targeted metaproteomic methods have been deployed in marine studies using stable isotope  
751 labeled peptides for calibration, achieving femtomoles per liter of seawater estimates of  
752 transporters, regulatory proteins, and enzymes (Saito et al., 2020; Bertrand et al., 2013; Saito et  
753 al., 2014, 2015; Joy-Warren et al., 2022; Wu et al., 2019). These methods are not yet widely  
754 adopted, but with growing interest could be deployed to other laboratories and incorporated into  
755 future iterations of intercomparison and intercalibration studies. DIA also has great potential in  
756 ocean metaproteome studies and is increasingly being deployed in laboratory and field studies  
757 of marine systems. Similar to this DDA intercomparison, the methodological and [bioinformatic](#)  
758 challenges of DIA could be explored during intercomparisons of analyses of ocean samples.  
759 Finally, as mentioned above, all participants of this study used orbitrap mass spectrometers for  
760 DDA submissions, but new instrumentation such as trapped ion mobility spectrometry time of  
761 flight mass spectrometers (timsTOF) may be applied to ocean metaproteome analyses and  
762 would be important to intercompare with orbitrap platforms.

763 As noted above, there were also challenges in collating and comparing data outputs  
764 from various software, as well as variation in how those programs conducted protein inference.  
765 For example, peptide-level data from different research groups were reported as either  
766 unmodified peptide sequences or as various peptide analytes (where modifications and charges  
767 states were included with the peptide sequence), making compilation of peptide reports difficult.  
768 Similarly, at the protein level reported proteins could be counted either before or after protein  
769 grouping, e.g. applying Occam's-razor logic to peptide groupings into proteins – the former

770 reflecting the set of all proteins in the database that could be in the sample, the latter the  
771 minimum set required to explain the peptide data. Such issues will also contribute to challenges  
772 in integration and assembly of data from different laboratories for large ocean datasets. While  
773 best practices for metadata and data types have been described by the community that include  
774 specific attributes important for environmental and ocean samples such as geospatial location  
775 and sample collection information (Saito et al., 2019) similar to the metadata standard recently  
776 put forward in the human proteome field (Dai et al., 2021), this study also demonstrated that  
777 there is ~~continues to be~~ a need for standardization of data output formats for metaproteomic  
778 results, ~~similar to the metadata standard recently put forward in the human proteome field (Dai~~  
779 et al., 2021).

#### 780 *4.5 Metaproteomics in Global Ocean Surveys*

781 Understanding how the oceans are responding to the rapid changes driven by human  
782 alteration of ecosystems is a high priority. Ocean and environmental sciences have a long  
783 history of chemical measurements that are critical to assessing ecosystems and climatic  
784 change. Such measurements have been straightforward for discrete measurements, such as  
785 temperature, pH, chlorophyll, phosphate, dissolved iron and numerous other variables. When  
786 collected over large spatial (ocean basin) or temporal (seasonal or decadal spans) scales, these  
787 datasets have been powerful in identifying major (both cyclical and secular) changes. 'Omics'  
788 measurements represent a more complex data type where each discrete sample can generate  
789 thousands (if not more) of units of information. This study demonstrates the power and potential  
790 for collaborative metaproteomics studies to identify key functional molecules and relate them to  
791 their taxonomic microbial sources within the microbiome from multiple lab groups. Moreover,  
792 multi-lab metaproteomics results in vastly enhanced identification of low abundance proteins  
793 that are not identified by all research groups. Such low abundance proteins can be more likely  
794 to change in abundance with changing environmental conditions and nutrient limitations,

795 resulting in a more nuanced and richer investigation of marine microbial ecology and  
796 biogeochemistry with collaborative metaproteomics research. The implementation of such  
797 voluminous data is beginning to be applied on larger scales and holds great promise in  
798 improving not only our understanding of the functioning of the current system, but also the way  
799 we assess how environments are changing with continued human perturbations.

800 Intercomparison and intercalibration are critical activities to undertake in order to allow  
801 comparison of omics results across time and space dimensions. With major programs  
802 underway and being envisioned such as the BioGEO TRACES, AtlantECO, Bio-GO-SHIP, and  
803 BioGeoSCAPES efforts, the imperative for such intercalibration has grown and the need for best  
804 practices is urgent. This Ocean Metaproteomic Intercomparison study is a valuable step in  
805 assessing metaproteomic capabilities across a number of international laboratories,  
806 demonstrating a clear consistency in measurement capability, while also pointing to the  
807 potential for continued community development of metaproteomic capacity and technology.

808

809 *Author Contributions*:- MAS and MRM obtained OCB workshop support and drafted the  
810 experimental design with feedback from BN, MJ, and DL acting as the Advisory Committee. SC,  
811 JH, DL, GJV, and JKS conducted the metagenomic analyses and assembly. JKS, MAS, MMB,  
812 MRM, and RM conducted data analysis and visualization. MRM, MAS, JAB, MVJ, and RJ  
813 conducted sample collection and/or AUV Clio operations. MAS, JKS, MRM, EMB, SC, JRC, TG,  
814 JH, RLH, PJ, MJ, RK, HK, DL, JSPM, EM, SM, DMM, JN, BN, JJ, MD, GJH, RG, RM, BLN, MP,  
815 SP, AR, ER, BS, TVDB, JRW, HZ, and ZZ contributed mass spectrometry data and/or  
816 bioinformatics data for the intercomparison. JKS anonymized data submissions and conducted  
817 follow-up correspondence about methods. The manuscript was drafted by MAS and all authors  
818 contributed to the writing and editing.

819



820 *Data and Code Availability*-. The raw files [\\_metagenome database](#)  
821 [\(Intercal\\_ORFs\\_prodigal\\_metagenemark.fasta\)](#), and associated annotations  
822 [\(Intercal\\_assembly\\_annotations.csv\)](#) for this project summarized in Table S3 are available at  
823 ProteomeXchange and PRIDE repository with the dataset identifier PXD043218 and  
824 [PXD04423440.6019/PXD043218](#). Access for reviewers is available using the username:  
825 [reviewer\\_pxd043218@ebi.ac.uk](#) and password: [uSxV/kRza](#), and  
826 [reviewer\\_pxd044234@ebi.ac.uk](#) and password: [Evvgeed0](#). Co-located information about these  
827 datasets are available at the Biological and Chemical Data Management Office under project  
828 765945 (<https://www.bco-dmo.org/project/765945>) and at the BATS page ([https://www.bco-](https://www.bco-dmo.org/project/2124)  
829 [dmo.org/project/2124](https://www.bco-dmo.org/project/2124)). The metagenomic reads are listed under Bioproject Accession:  
830 PRJNA932835; SRA submission: SUB12819843, available at link:  
831 <https://www.ncbi.nlm.nih.gov/bioproject/PRJNA932835>. The code for ~~fer~~ upset visualization is  
832 available at: <https://maggimars.github.io/protein/PeptideUpSetR.html>.

833

834 *Competing Interests* - The authors declare no competing financial interests.

835 *Supplemental Materials* - Methods for the [bio](#)informatic intercomparison study are available in  
836 the Supplemental Methods. Supplemental Information is available as Tables S1-S11, and  
837 Figures S1-S3.

838 *Acknowledgements* - This manuscript is a product of the sustained efforts of a small group  
839 activity supported by the Ocean Carbon & Biogeochemistry (OCB) Project Office (NSF OCE-  
840 1850983 and NASA NNX17AB17G), based on a proposal written by M.A.S. and M.R.M. The  
841 research expedition where samples were collected was supported by the NSF Biological  
842 Oceanography and Chemical Oceanography. We also thank the R/V *Atlantic Explorer* and the  
843 Bermuda Atlantic Time-series Study team for assistance at sea. AUV Clio sample collection was  
844 supported by NSF OCE 1658030 and 1924554. Analyses by participating laboratories

845 acknowledge support from: NSERC Discovery Grant RGPIN-2015-05009 and Simons  
846 Foundation Grant 504183 to E.M.B, the Austrian Science Fund (FWF) DEPOCA (project  
847 number AP3558721) to G.J.H., Simons Foundation grant 402971 to J.R.W., National Institute of  
848 Health 1R21ES034337-01 to B.L.N., the Norwegian Centennial Chair Program at the University  
849 of Minnesota for funding to PDJ, SM, and TJG, NIH R01 GM135709, NSF OCE-1924554, OCE-  
850 2019589 and Simons Foundation Grant 1038971 to M.A.S. Identification of certain commercial  
851 equipment, instruments, software, or materials does not imply recommendation or endorsement  
852 by the National Institute of Standards and Technology, nor does it imply that the products  
853 identified are necessarily the best available for the purpose. We thank Magnus Palmblad, John  
854 Kucklick, and an anonymous reviewer for comments on the [pre-submission version of the](#)  
855 manuscript. [We also thank two anonymous reviewers for their constructive comments during](#)  
856 [manuscript review.](#)

857

858

859

860 **References**

- 861 Bender, S. J., Moran, D. M., McIlvin, M. R., Zheng, H., McCrow, J. P., Badger, J., DiTullio, G.  
862 R., Allen, A. E., and Saito, M. A.: Colony formation in *Phaeocystis antarctica*: connecting  
863 molecular mechanisms with iron biogeochemistry, *Biogeosciences*, 15, 4923–4942, 2018.
- 864 Bergauer, K., Fernandez-Guerra, A., Garcia, J. A., Sprenger, R. R., Stepanauskas, R.,  
865 Pachiadaki, M. G., Jensen, O. N., and Herndl, G. J.: Organic matter processing by microbial  
866 communities throughout the Atlantic water column as revealed by metaproteomics, *Proceedings*  
867 *of the National Academy of Sciences*, 115, E400–E408, 2018.
- 868 Bertrand, E. M., Moran, D. M., McIlvin, M. R., Hoffman, J. M., Allen, A. E., and Saito, M. A.:  
869 Methionine synthase interreplacement in diatom cultures and communities: Implications for the  
870 persistence of B12 use by eukaryotic phytoplankton, *Limnology and Oceanography*, 58, 1431–  
871 1450, 2013.
- 872 Breier, J. A., Jakuba, M. V., Saito, M. A., Dick, G. J., Grim, S. L., Chan, E. W., McIlvin, M. R.,  
873 Moran, D. M., Alanis, B. A., and Allen, A. E.: Revealing ocean-scale biochemical structure with a  
874 deep-diving vertical profiling autonomous vehicle, *Science Robotics*, 5, eabc7104, 2020.
- 875 Buchfink, B., Xie, C., and Huson, D. H.: Fast and sensitive protein alignment using DIAMOND,  
876 *Nature methods*, 12, 59–60, 2015.
- 877 [Carlson, C.A., Morris, R., Parsons, R., Treusch, A.H., Giovannoni, S.J. and Vergin, K., 2009.](#)  
878 [Seasonal dynamics of SAR11 populations in the euphotic and mesopelagic zones of the](#)  
879 [northwestern Sargasso Sea. \*The ISME journal\*, 3\(3\), pp.283-295.](#)
- 880 [Casey, J.R., Lomas, M.W., Mandecki, J. and Walker, D.E., 2007. \*Prochlorococcus contributes\*](#)  
881 [to new production in the Sargasso Sea deep chlorophyll maximum. \*Geophysical Research\*](#)  
882 [Letters](#), 34(10).
- 883
- 884 Cohen, N. R., McIlvin, M. R., Moran, D. M., Held, N. A., Saunders, J. K., Hawco, N. J.,  
885 Brosnahan, M., DiTullio, G. R., Lamborg, C., and McCrow, J. P.: Dinoflagellates alter their  
886 carbon and nutrient metabolic strategies across environmental gradients in the central Pacific  
887 Ocean, *Nature Microbiology*, 6, 173–186, 2021.
- 888 Cohen, N. R., Krinos, A. I., Kell, R. M., Chmiel, R. J., Moran, D. M., McIlvin, M. R., Lopez, P. Z.,  
889 Barth, A., Stone, J., Alanis, B. A., Chan, E. W., Breier, J. A., Jakuba, M. V., Johnson, R.,  
890 Alexander, H., and Saito, M. A.: Microeukaryote metabolism across the western North Atlantic  
891 Ocean revealed through autonomous underwater profiling, *Ecology*,  
892 <https://doi.org/10.1101/2023.11.20.567900>, 2023.
- 893 Coleman, M. L. and Chisholm, S. W.: Ecosystem-specific selection pressures revealed through  
894 comparative population genomics, *Proceedings of the National Academy of Sciences*, 107,  
895 18634–18639, 2010.
- 896 Conway, J. R., Lex, A., and Gehlenborg, N.: UpSetR: an R package for the visualization of  
897 intersecting sets and their properties, *Bioinformatics*, 2017.

Formatted: Normal

898 Dai, C., Füllgrabe, A., Pfeuffer, J., Solovyeva, E. M., Deng, J., Moreno, P., Kamatchinathan, S.,  
899 Kundu, D. J., George, N., and Fexova, S.: A proteomics sample metadata representation for  
900 multiomics integration and big data analysis, *Nature Communications*, 12, 1–8, 2021.

901 Deutsch, E. W., Bandeira, N., Sharma, V., Perez-Riverol, Y., Carver, J. J., Kundu, D. J., García-  
902 Seisdedos, D., Jarnuczak, A. F., Hewapathirana, S., Pullman, B. S., Wertz, J., Sun, Z., Kawano,  
903 S., Okuda, S., Watanabe, Y., Hermjakob, H., MacLean, B., MacCoss, M. J., Zhu, Y., Ishihama,  
904 Y., and Vizcaíno, J. A.: The ProteomeXchange consortium in 2020: enabling ‘big data’  
905 approaches in proteomics, *Nucleic Acids Research*, gkz984, <https://doi.org/10.1093/nar/gkz984>,  
906 2019.

907 [Falkowski, P.G., Fenchel, T. and Delong, E.F., 2008. The microbial engines that drive Earth's  
908 biogeochemical cycles. \*science\*, 320\(5879\), 1034-1039.](#)

909

910 Fuchsman, C. A., Palevsky, H. I., Widner, B., Duffy, M., Carlson, M. C., Neibauer, J. A.,  
911 Mulholland, M. R., Keil, R. G., Devol, A. H., and Rocap, G.: Cyanobacteria and cyanophage  
912 contributions to carbon and nitrogen cycling in an oligotrophic oxygen-deficient zone, *The ISME  
913 Journal*, 13, 2714–2726, 2019.

914 Georges, A. A., El-Swais, H., Craig, S. E., Li, W. K., and Walsh, D. A.: Metaproteomic analysis  
915 of a winter to spring succession in coastal northwest Atlantic Ocean microbial plankton, *The  
916 ISME journal*, 8, 1301–1313, 2014.

917 Hawley, A. K., Brewer, H. M., Norbeck, A. D., Paša-Tolić, L., and Hallam, S. J.: Metaproteomics  
918 reveals differential modes of metabolic coupling among ubiquitous oxygen minimum zone  
919 microbes, *Proceedings of the National Academy of Sciences*, 111, 11395–11400, 2014.

920 Held, N. A., Sutherland, K. M., Webb, E. A., McIlvin, M. R., Cohen, N. R., Devaux, A. J.,  
921 Hutchins, D. A., Waterbury, J. B., Hansel, C. M., and Saito, M. A.: Mechanisms and  
922 heterogeneity of in situ mineral processing by the marine nitrogen fixer *Trichodesmium* revealed  
923 by single-colony metaproteomics, *ISME Communications*, 1, 1–9, 2021.

924 [Hulstaert, N., Shofstahl, J., Sachsenberg, T., Walzer, M., Barsnes, H., Martens, L. and Perez-  
925 Riverol, Y., 2019. ThermoRawFileParser: modular, scalable, and cross-platform RAW file  
926 conversion. \*Journal of Proteome Research\*, 19\(1\), 537-542.](#)

927

928 Hyatt, D., Chen, G.-L., LoCascio, P. F., Land, M. L., Larimer, F. W., and Hauser, L. J.: Prodigal:  
929 prokaryotic gene recognition and translation initiation site identification, *BMC bioinformatics*, 11,  
930 1–11, 2010.

931 Jagtap, P. D., Hoopmann, M. R., Neely, B. A., Harvey, A., Käll, L., Perez-Riverol, Y., Abajorga,  
932 M. K., Thomas, J. A., Weintraub, S. T., and Palmblad, M.: The Association of Biomolecular  
933 Resource Facilities Proteome Informatics Research Group Study on Metaproteomics (iPRG-  
934 2020), *J Biomol Tech*, 34, 3fc1f5fe.a058bad4, <https://doi.org/10.7171/3fc1f5fe.a058bad4>, 2023.

935 Joy-Warren, H. L., Alderkamp, A.-C., van Dijken, G. L., J Jabre, L., Bertrand, E. M., Baldonado,  
936 E. N., Glickman, M. W., Lewis, K. M., Middag, R., and Seyitmuhammedov, K.: Springtime  
937 phytoplankton responses to light and iron availability along the western Antarctic Peninsula,  
938 *Limnology and Oceanography*, 67, 800–815, 2022.

Formatted: Normal

Formatted: Font: 11 pt

Formatted: Normal

939 Kanehisa, M., Sato, Y., and Morishima, K.: BlastKOALA and GhostKOALA: KEGG tools for  
940 functional characterization of genome and metagenome sequences, *Journal of molecular*  
941 *biology*, 428, 726–731, 2016.

942  
943 [Keller, A., Nesvizhskii, A.I., Kolker, E. and Aebersold, R., 2002. An explanation of the Peptide  
944 Prophet algorithm developed. \*Anal. Chem.\*, 74\(2002\), 5383-5392.](#)

945  
946 [Kim, S. and Pevzner, P.A., 2014. MS-GF+ makes progress towards a universal database  
947 search tool for proteomics. \*Nature Communications\*, 5\(1\), 5277.](#)

948  
949 [Kiweler, M., Looso, M. and Graumann, J., 2019. MARMoSET—extracting publication-ready mass  
950 spectrometry metadata from RAW files. \*Molecular & Cellular Proteomics\*, 18\(8\), 1700-1702.](#)

951  
952 Kleiner, M., Thorson, E., Sharp, C. E., Dong, X., Liu, D., Li, C., and Strous, M.: Assessing  
953 species biomass contributions in microbial communities via metaproteomics, *Nature*  
954 *Communications*, 8, 1–14, 2017.

955  
956 [Kleiner, M., 2019. Metaproteomics: much more than measuring gene expression in microbial  
957 communities. \*Msystems\*, 4\(3\), 1128/msystems.00115-19.](#)

958 Leary, D. H., Li, R. W., Hamdan, L. J., Hervey IV, W. J., Lebedev, N., Wang, Z., Deschamps, J.  
959 R., Kusterbeck, A. W., and Vora, G. J.: Integrated metagenomic and metaproteomic analyses of  
960 marine biofilm communities, *Biofouling*, 30, 1211–1223, 2014.

961  
962 Malmstrom, R. R., Coe, A., Kettler, G. C., Martiny, A. C., Frias-Lopez, J., Zinser, E. R., and  
963 Chisholm, S. W.: Temporal dynamics of *Prochlorococcus* ecotypes in the Atlantic and Pacific  
oceans, *The ISME journal*, 4, 1252–1264, 2010.

964  
965 McCain, J. S. P. and Bertrand, E. M.: Prediction and consequences of cofragmentation in  
metaproteomics, *Journal of proteome research*, 18, 3555–3566, 2019.

966  
967 McCain, J. S. P., Allen, A. E., and Bertrand, E. M.: Proteomic traits vary across taxa in a coastal  
Antarctic phytoplankton bloom, *The ISME journal*, 16, 569–579, 2022.

968  
969 McIlvin, M. R. and Saito, M. A.: Online Nanoflow Two-Dimension Comprehensive Active  
970 Modulation Reversed Phase–Reversed Phase Liquid Chromatography High-Resolution Mass  
971 Spectrometry for Metaproteomics of Environmental and Microbiome Samples, *Journal of*  
*proteome research*, 20, 4589–4597, 2021.

972  
973 Mikan, M. P., Harvey, H. R., Timmins-Schiffman, E., Riffle, M., May, D. H., Salter, I., Noble, W.  
974 S., and Nunn, B. L.: Metaproteomics reveal that rapid perturbations in organic matter prioritize  
975 functional restructuring over taxonomy in western Arctic Ocean microbiomes, *The ISME journal*,  
14, 39–52, 2020.

976  
977 Moore, E. K., Nunn, B. L., Goodlett, D. R., and Harvey, H. R.: Identifying and tracking proteins  
978 through the marine water column: Insights into the inputs and preservation mechanisms of  
protein in sediments, *Geochimica et cosmochimica acta*, 83, 324–359, 2012.

Formatted: Font: 11 pt

Formatted: Font: 11 pt

Formatted: Normal

Formatted: Font: 11 pt

Formatted: Font: Arial, 11 pt

Formatted: Font: 11 pt

Formatted: Normal

979  
980  
981  
982  
983  
984  
985  
986  
  
987  
988  
  
989  
990  
991  
992  
993  
994  
995  
  
996  
997  
998  
  
999  
1000  
1001  
  
1002  
1003  
1004  
  
1005  
1006  
1007  
  
1008  
1009  
1010  
1011  
  
1012  
1013  
1014  
  
1015  
1016  
1017  
1018  
1019

[Moran, M.A., Kujawinski, E.B., Schroer, W.F., Amin, S.A., Bates, N.R., Bertrand, E.M., Braakman, R., Brown, C.T., Covert, M.W., Doney, S.C. and Dyhrman, S.T., 2022. Microbial metabolites in the marine carbon cycle. \*Nature microbiology\*, 7\(4\), 508-523.](#)

Morris, R. M., Nunn, B. L., Frazar, C., Goodlett, D. R., Ting, Y. S., and Rocap, G.: Comparative metaproteomics reveals ocean-scale shifts in microbial nutrient utilization and energy transduction, *The ISME journal*, 4, 673–685, 2010.

Mueller, R. S. and Pan, C.: Sample handling and mass spectrometry for microbial metaproteomic analyses, in: *Methods in Enzymology*, vol. 531, Elsevier, 289–303, 2013.

[Nesvizhskii, A.I., Keller, A., Kolker, E. and Aebersold, R., 2003. A statistical model for identifying proteins by tandem mass spectrometry. \*Analytical Chemistry\*, 75\(17\), 4646-4658.](#)

Piehowski, P. D., Petyuk, V. A., Orton, D. J., Xie, F., Moore, R. J., Ramirez-Restrepo, M., Engel, A., Lieberman, A. P., Albin, R. L., and Camp, D. G.: Sources of technical variability in quantitative LC–MS proteomics: human brain tissue sample analysis, *Journal of proteome research*, 12, 2128–2137, 2013.

Pietilä, S., Suomi, T., and Elo, L. L.: Introducing untargeted data-independent acquisition for metaproteomics of complex microbial samples, *ISME COMMUN.*, 2, 51, <https://doi.org/10.1038/s43705-022-00137-0>, 2022.

Ram, R. J., VerBerkmoes, N. C., Thelen, M. P., Tyson, G. W., Baker, B. J., Blake, R. C., Shah, M., Hettich, R. L., and Banfield, J. F.: Community proteomics of a natural microbial biofilm, *Science*, 308, 1915–1920, 2005.

Saito, M. A., McIlvin, M. R., Moran, D. M., Goepfert, T. J., DiTullio, G. R., Post, A. F., and Lamborg, C. H.: Multiple nutrient stresses at intersecting Pacific Ocean biomes detected by protein biomarkers, *Science*, 345, 1173–1177, 2014.

Saito, M. A., Dorsk, A., Post, A. F., McIlvin, M. R., Rappé, M. S., DiTullio, G. R., and Moran, D. M.: Needles in the blue sea: Sub-species specificity in targeted protein biomarker analyses within the vast oceanic microbial metaproteome, *Proteomics*, 15, 3521–3531, 2015.

Saito, M. A., Bertrand, E. M., Duffy, M. E., Gaylord, D. A., Held, N. A., Hervey IV, W. J., Hettich, R. L., Jagtap, P. D., Janech, M. G., and Kinkade, D. B.: Progress and challenges in ocean metaproteomics and proposed best practices for data sharing, *Journal of proteome research*, 18, 1461–1476, 2019.

Saito, M. A., McIlvin, M. R., Moran, D. M., Santoro, A. E., Dupont, C. L., Rafter, P. A., Saunders, J. K., Kaul, D., Lamborg, C. H., and Westley, M.: Abundant nitrite-oxidizing metalloenzymes in the mesopelagic zone of the tropical Pacific Ocean, *Nature Geoscience*, 13, 355–362, 2020.

Saunders, J. K., Gaylord, D. A., Held, N. A., Symmonds, N., Dupont, C. L., Shepherd, A., Kinkade, D. B., and Saito, M. A.: METATRYP v 2.0: Metaproteomic least common ancestor analysis for taxonomic inference using specialized sequence assemblies—standalone software and web servers for marine microorganisms and coronaviruses, *Journal of proteome research*, 19, 4718–4729, 2020.

Formatted: Font: 11 pt

Formatted: Normal

Formatted: Font: 11 pt

Formatted: Normal

1020 Scanlan, D. J., Silman, N. J., Donald, K. M., Wilson, W. H., Carr, N. G., Joint, I., and Mann, N.  
1021 H.: An immunological approach to detect phosphate stress in populations and single cells of  
1022 photosynthetic picoplankton, *Applied and environmental microbiology*, 63, 2411–2420, 1997.

1023 Schiebenhoefer, H., Van Den Bossche, T., Fuchs, S., Renard, B. Y., Muth, T., and Martens, L.:  
1024 Challenges and promise at the interface of metaproteomics and genomics: an overview of  
1025 recent progress in metaproteogenomic data analysis, *Expert Review of Proteomics*, 16, 375–  
1026 390, 2019.

1027 Sørensen, T.: A method of establishing groups of equal amplitude in plant sociology based on  
1028 similarity of species and its application to analyses of the vegetation on Danish common.,  
1029 *Kongelige Danske Videnskabernes Selskab*, 5, 1–34, 1948.

1030 Sowell, S. M., Wilhelm, L. J., Norbeck, A. D., Lipton, M. S., Nicora, C. D., Barofsky, D. F.,  
1031 Carlson, C. A., Smith, R. D., and Giovanonni, S. J.: Transport functions dominate the SAR11  
1032 metaproteome at low-nutrient extremes in the Sargasso Sea, *The ISME journal*, 3, 93–105,  
1033 2009.

1034 Stewart, H. I., Grinfeld, D., Giannakopoulos, A., Petzoldt, J., Shanley, T., Garland, M., Denisov,  
1035 E., Peterson, A. C., Damoc, E., Zeller, M., Arrey, T. N., Pashkova, A., Renuse, S., Hakimi, A.,  
1036 Kühn, A., Biel, M., Kreuzmann, A., Hagedorn, B., Colonius, I., Schütz, A., Stefes, A., Dwivedi,  
1037 A., Mourad, D., Hoek, M., Reitemeier, B., Cochems, P., Kholomeev, A., Ostermann, R., Quiring,  
1038 G., Ochmann, M., Möhring, S., Wagner, A., Petker, A., Kanngiesser, S., Wiedemeyer, M.,  
1039 Balschun, W., Hermanson, D., Zabrouskov, V., Makarov, A. A., and Hock, C.: Parallelized  
1040 Acquisition of Orbitrap and Astral Analyzers Enables High-Throughput Quantitative Analysis,  
1041 *Anal. Chem.*, 95, 15656–15664, <https://doi.org/10.1021/acs.analchem.3c02856>, 2023.

1042 Tagliabue, A.: ‘Oceans are hugely complex’: modelling marine microbes is key to climate  
1043 forecasts, *Nature*, 623, 250–252, <https://doi.org/10.1038/d41586-023-03425-4>, 2023.

1044 Timmins-Schiffman, E., May, D. H., Mikan, M., Riffle, M., Frazar, C., Harvey, H. R., Noble, W.  
1045 S., and Nunn, B. L.: Critical decisions in metaproteomics: achieving high confidence protein  
1046 annotations in a sea of unknowns, *The ISME journal*, 11, 309–314, 2017.

1047 Ustick, L. J., Larkin, A. A., Garcia, C. A., Garcia, N. S., Brock, M. L., Lee, J. A., Wiseman, N. A.,  
1048 Moore, J. K., and Martiny, A. C.: Metagenomic analysis reveals global-scale patterns of ocean  
1049 nutrient limitation, *Science*, 372, 287–291, 2021.

1050 Van Den Bossche, T., Kunath, B. J., Schallert, K., Schäpe, S. S., Abraham, P. E., Armengaud,  
1051 J., Arntzen, M. Ø., Bassignani, A., Benndorf, D., and Fuchs, S.: Critical Assessment of  
1052 MetaProteome Investigation (CAMPI): a multi-laboratory comparison of established workflows,  
1053 *Nature communications*, 12, 1–15, 2021.

1054 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D.,  
1055 Burovski, E., Peterson, P., Weckesser, W., and Bright, J.: SciPy 1.0: fundamental algorithms for  
1056 scientific computing in Python, *Nature methods*, 17, 261–272, 2020.

1057 Waskom, M. L.: Seaborn: statistical data visualization, *Journal of Open Source Software*, 6,  
1058 3021, 2021.

1059 Williams, T. J., Long, E., Evans, F., DeMaere, M. Z., Lauro, F. M., Raftery, M. J., Ducklow, H.,  
1060 Grzymski, J. J., Murray, A. E., and Cavicchioli, R.: A metaproteomic assessment of winter and  
1061 summer bacterioplankton from Antarctic Peninsula coastal surface waters, *The ISME journal*, 6,  
1062 1883–1900, 2012.

1063 Wilmes, P. and Bond, P. L.: Metaproteomics: studying functional gene expression in microbial  
1064 ecosystems, *Trends in microbiology*, 14, 92–97, 2006.

1065 Wilmes, P., Andersson, A. F., Lefsrud, M. G., Wexler, M., Shah, M., Zhang, B., Hettich, R. L.,  
1066 Bond, P. L., VerBerkmoes, N. C., and Banfield, J. F.: Community proteogenomics highlights  
1067 microbial strain-variant protein expression within activated sludge performing enhanced  
1068 biological phosphorus removal, *The ISME journal*, 2, 853–864, 2008.

1069 [Worden, A.Z., Follows, M.J., Giovannoni, S.J., Wilken, S., Zimmerman, A.E. and Keeling, P.J.,](#)  
1070 [2015. Rethinking the marine carbon cycle: factoring in the multifarious lifestyles of](#)  
1071 [microbes. \*Science\*, 347\(6223\), 1257594.](#)

1072  
1073 Wu, M., McCain, J. S. P., Rowland, E., Middag, R., Sandgren, M., Allen, A. E., and Bertrand, E.  
1074 M.: Manganese and iron deficiency in Southern Ocean *Phaeocystis antarctica* populations  
1075 revealed through taxon-specific protein indicators, *Nature communications*, 10, 1–10, 2019.

1076  
1077

Formatted: Font: 11 pt

Formatted: Normal



1078 **Figure Captions**

1079 **Figure 1.** Ocean metaproteomics intercomparison experimental design and sample collection.

1080 a) The laboratory component (left) consisted of collection of field samples, 1-dimensional (1D)  
1081 chromatographic separation followed by data dependent analysis (DDA) uniformly employing  
1082 orbitrap mass spectrometers analyses by participating laboratories and submission of raw and  
1083 processed data. The [bioinformatic](#) (right) component consisted of distribution of two 1D-DDA  
1084 files, peptide-to-spectrum matching (PSMs), and submission and compilation of results. b) Size-  
1085 fractionated sample collection on 3.0  $\mu\text{m}$  pore-size filter followed by a 0.2  $\mu\text{m}$  pore-size Supor  
1086 filter, and the 0.2–3.0  $\mu\text{m}$  size fraction was used for the intercomparison study. c) Two horizontal  
1087 *in-situ* McLane pumps were bracketed together with two Mini-MULVS filter head units each and  
1088 deployment on synthetic line. d) The four 142 mm filters were sliced into eighths (inset) and two  
1089 slices were distributed to each participating laboratory.

1090

1091 **Figure 2.** Shared peptides and proteins between laboratory groups using laboratory

1092 submissions processed through a single [bioinformatics](#) re-analysis pipeline. a) Total number of  
1093 discovered unique peptides varied by more than three-fold among seven laboratory groups  
1094 (horizontal bars) due to varying extraction and analytical schemes (FDR 0.1%). The number of  
1095 intersections between datasets across all seven datasets was 1395 (fourth blue bar from left),  
1096 and various sets of intersections of peptides were observed amongst the data. b) Total number  
1097 of discovered proteins (FDR < 1%) varied more than four-fold from 1586 to 6221 among labs  
1098 (horizontal bars). Intersections between datasets across all seven laboratories was 1056, with  
1099 various sets of intersections of proteins observed, similar to the peptides. c) 7-way Venn  
1100 diagrams of shared unique peptides between laboratories showed 1056 shared peptides  
1101 between the 7 laboratories. d) 3-way Venn diagrams showed 2398, 2304, and 3016 shared  
1102 unique peptides between laboratories.

1103

1104 **Figure 3.** Comparison of unique peptides and discovered proteins. Comparison as total protein  
1105 identifications and protein groups from the single pipeline re-analysis based on submissions  
1106 from 9 laboratories. Increasing sample depth is linear with mapping to proteins, ( $R^2$  of 0.97 and  
1107 0.98 for total protein IDs and protein groups, respectively, with slopes of 0.37 and 33) implying  
1108 that additional peptide discovery leads to proportionally more protein discovery, and that protein  
1109 discovery has not yet begun to saturate with more peptides mapping to each protein. Because  
1110 simple 1D analyses were stipulated in the intercomparison experimental design, peptide and  
1111 protein discovery was correspondingly limited in depth.

1112

1113 **Figure 4.** Quantitative comparison of intercomparison results. a) Pairwise comparisons of  
1114 ~~quantitative~~quantitative abundance across six laboratories in units of spectral counts  
1115 (comparisons with itself show unison diagonals). b)  $R^2$  values from pairwise linear regressions.  
1116 d) Total proteins identified in each laboratory. d) Average of each laboratory's  $R^2$  values from  
1117 pairwise regression with the other six laboratories (error bars are standard deviation). In all  
1118 cases average  $R^2$  values are higher than 0.5. e) Occurrences of  $R^2$  values in pairwise  
1119 comparisons spanning 0.4 to 0.9. Potential causes of this range are outlined in the Discussion  
1120 section.

1121

1122 **Figure 5.** Taxonomic and functional analysis of metaproteomic intercomparison. a) Percent  
1123 spectral counts by taxonomy was similar across laboratories and technical replicates within  
1124 laboratories. The sample was dominated by cyanobacteria and alphaproteobacteria,  
1125 corresponding primarily to *Prochlorococcus* and *Pelagibacter*, respectively. b) Percent spectral  
1126 counts per Kegg Ontology group showed the functional diversity of the sample.

1127

1128 **Figure 6.** Quantitative Sørensen similarity analysis. Analysis of top 1000 proteins (~75% of all

Formatted: Superscript

1129 proteins) showed 70–80% similarity between most laboratory groups. Technical triplicates for  
1130 each laboratory group are shown.

1131

1132 **Figure 7.** Intercomparison of [bio](#)informatic pipelines among laboratories. Unique peptide  
1133 identifications for sample Ocean 8 from 120m depth (a) and Ocean 11 from 20m depth (b), both  
1134 from the North Atlantic Ocean (Table S3), using a variety of pipelines and PSM algorithms.

1135

1136 **Figure 8.** Variability in discovered proteins between laboratories occurs in lower abundance  
1137 proteins. Top 7 panels: Abundance of proteins as percentage of total protein spectral counts  
1138 within each laboratory (y-axis is percentage), with proteins on the x-axis shown by ranked  
1139 abundance as the sum of spectral counts across all laboratories. Almost all proteins fall below  
1140 1% of spectral counts within the sample, and deeper proteomes have lower percentages due to  
1141 sharing of percent spectral counts across more discovered proteins. Bottom panel: Shared  
1142 proteins were found early within the long-tail of discovered proteins: the 1056 proteins shared  
1143 between all laboratory groups are almost all found to the left side indicating their higher  
1144 abundance in all seven datasets. Scale is binary in the seventh panel indicating presence in 7  
1145 labs or not.

1146

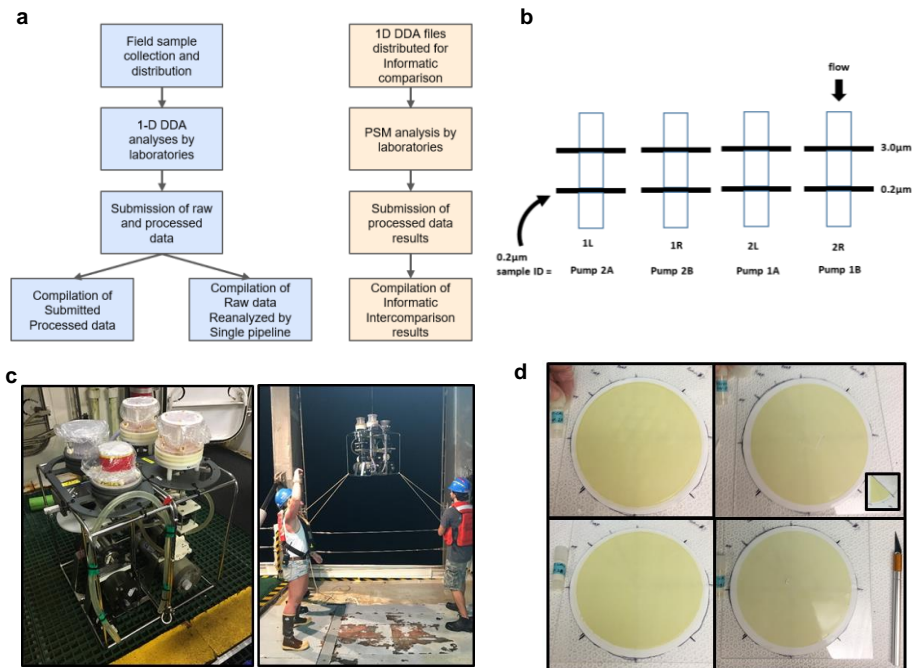
1147

1148 Figure 1.

1149

1150

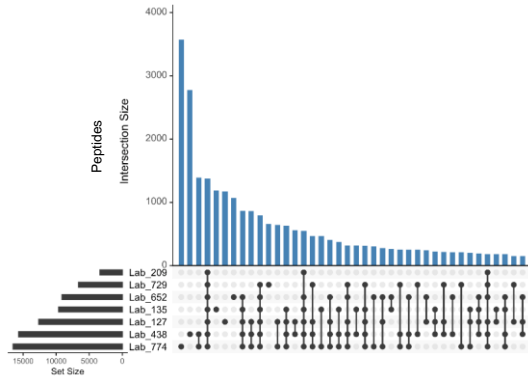
1151



1152 Figure 2.

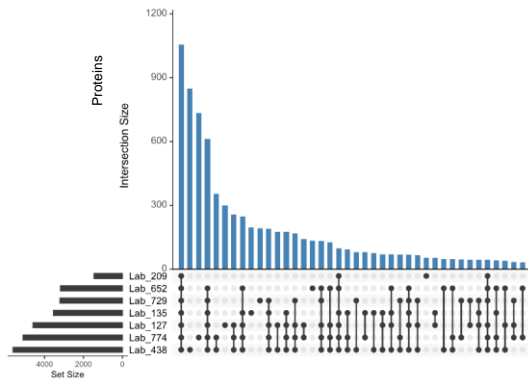
1153

1154



1156

1157

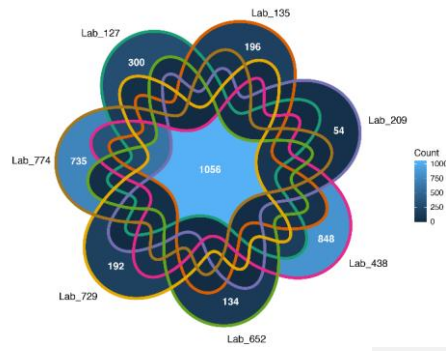


1163

1164

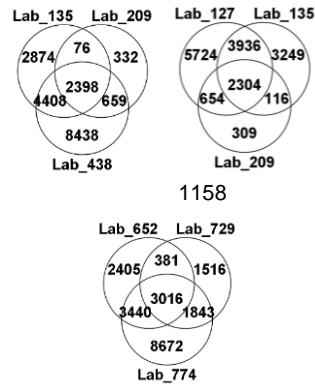
1165

c



1155

d



1166 Figure 3

1167

1168

1169

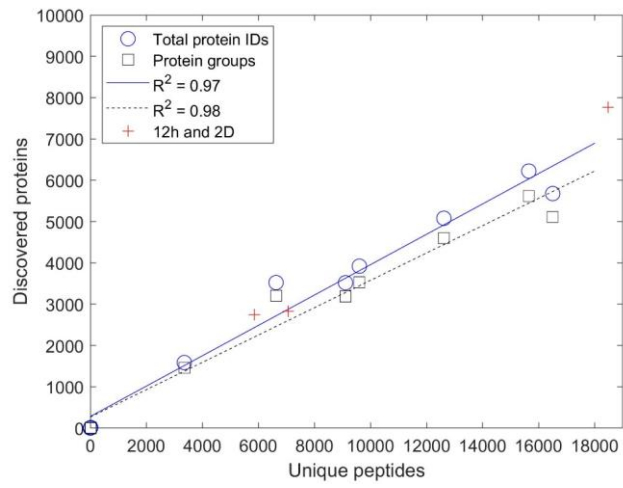
1170

1171

1172

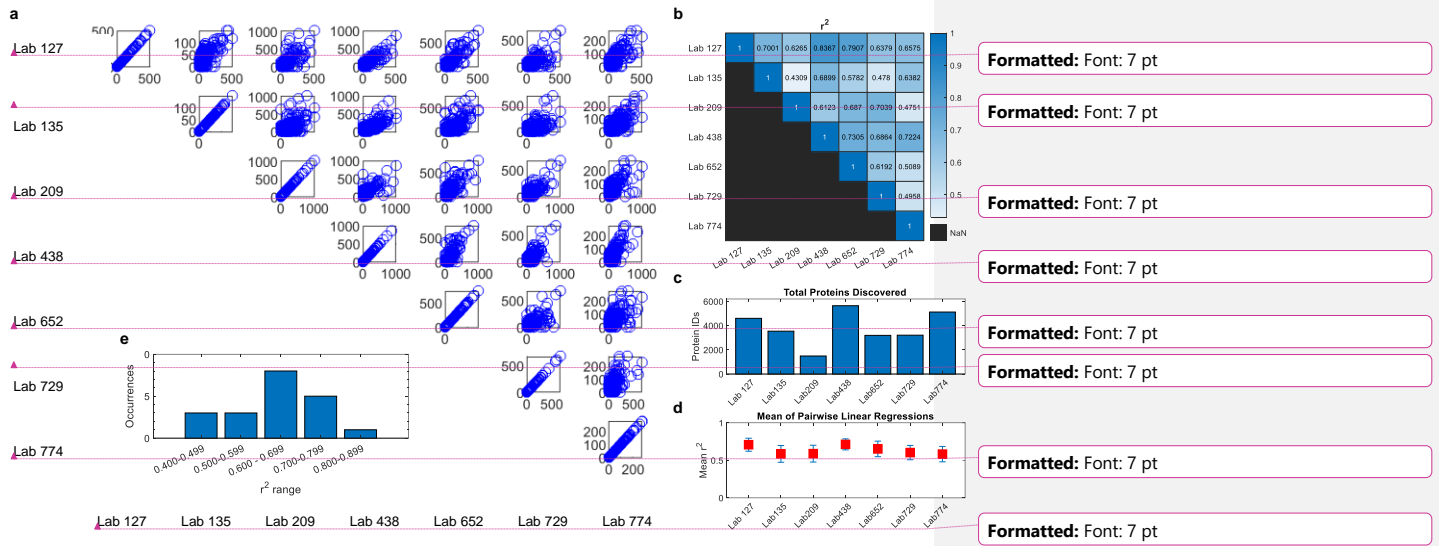
1173

1174



1175 Figure 4.

1176



Formatted: Font: 7 pt

Formatted: Font: 7 pt

Formatted: Font: 7 pt

Formatted: Font: 7 pt

Formatted: Font: 7 pt

Formatted: Font: 7 pt

Formatted: Font: 7 pt

Formatted: Font: 7 pt

1177

1178

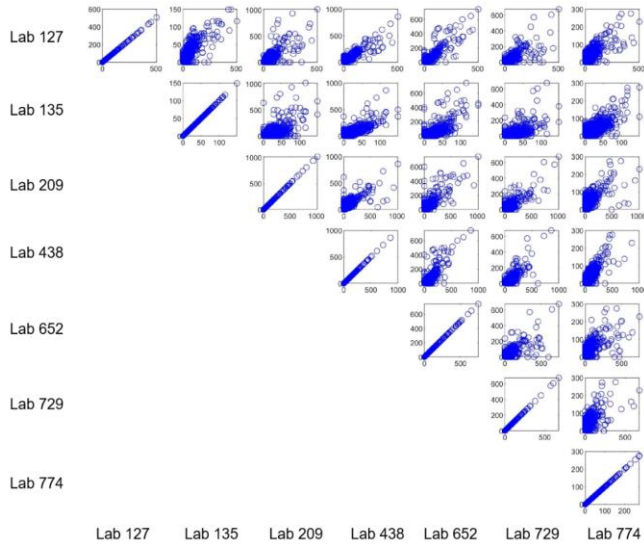
1179

1180

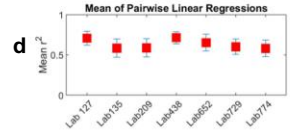
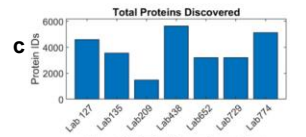
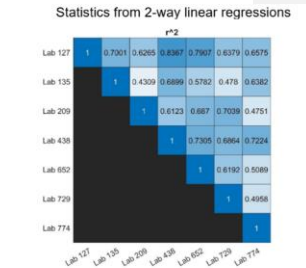
1181

1182

**a**



**b**





1183 Figure 5.

1184

1185

1186

1187

1188

1189

1190

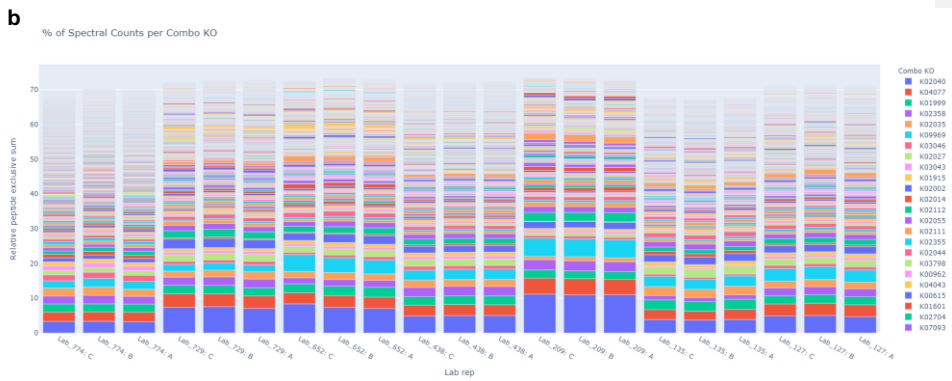
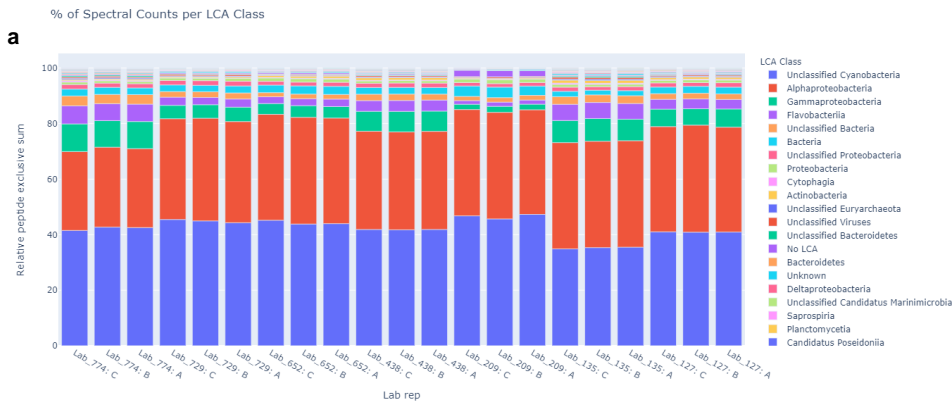
1191

1192

1193

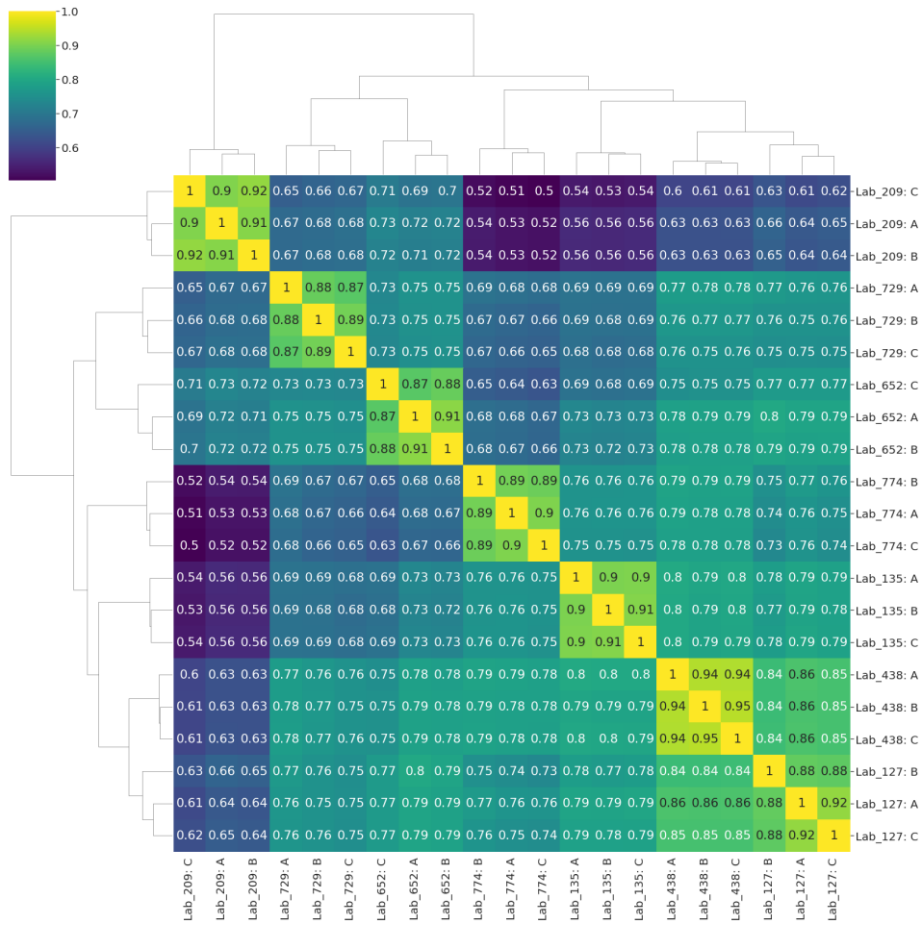
1194

1195



1196 Figure 6.

1197



1198

1199

1200

1201

1202

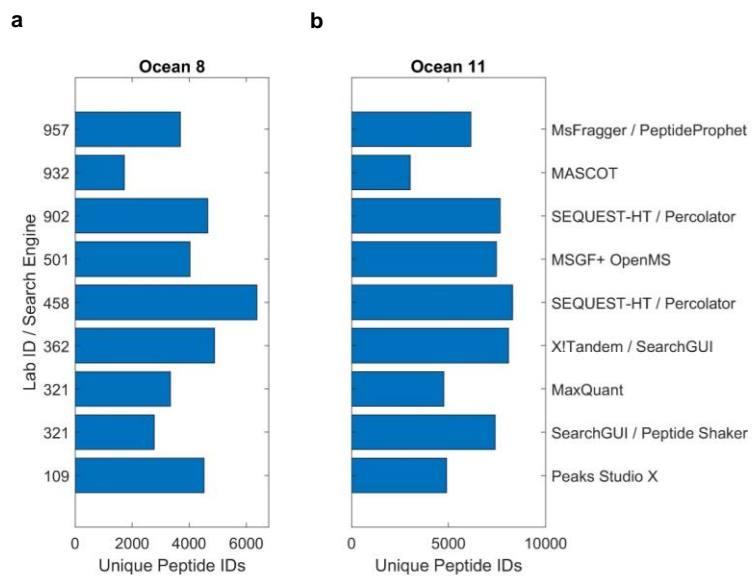
1203

1204 Figure 7.

1205

1206

1207



1208

1209

1210 Figure 8.

1211

