

Reviews

pubs.acs.org/jpr

¹ Progress and Challenges in Ocean Metaproteomics and Proposed 2 Best Practices for Data Sharing

- 3 Mak A. Saito,*,[†],[†] Erin M. Bertrand,[‡] Megan E. Duffy,[§] David A. Gaylord,[†] Noelle A. Held,[†] William Judson Hervey, IV, Bobert L. Hettich, Pratik Jagtap, Michael G. Janech, Danie B. Kinkade, Dasha Leary, Matthew McIlvin, Eli Moore, Robert Morris, Benjamin A. Neely, Brook Nunn, Jaclyn K. Saunders, Adam Shepherd, Adam Shepherd, Dasha Leary, Adam Shepherd, Dasha Leary, Daclyn K. Saunders, Adam Shepherd, Dasha Leary, Daclyn K. Saunders, Adam Shepherd, Dasha Leary, Daclyn K. Saunders, Dasha Leary, Dasha

- 7 Nicholas Symmonds, [†] and David Walsh
- 8 [†]Woods Hole Oceanographic Institution, Woods Hole, Massachusetts 02543, United States
- 9 [‡]Department of Biology, Dalhousie University, Halifax, Nova Scotia B3H 4R2, Canada
- 10 School of Oceanography, University of Washington, Seattle, Washington 98195-7940, United States
- 11 U.S. Naval Research Laboratory, Washington, D.C. 20375, United States
- ¹Oak Ridge National Laboratory and Microbiology Department, University of Tennessee, Knoxville, Tennessee 37996 United States
- 13 *Department of Biochemistry, Molecular Biology and Biophysics, University of Minnesota, Saint Paul, Minnesota 55108, United

20

21

22

23

24

2.5

26

2.7

28

29

30

31 32

33

34

35

- $^{
 abla}$ College of Charleston, Charleston, South Carolina 29424, United States
- Opepartment of Environmental Science, Rowan University, Glassboro, New Jersey 08028, United States
- National Institute of Standards and Technology, Gaithersburg, Maryland 20899, United States
- Department of Genome Sciences, University of Washington, Seattle, Washington 98195, United States
- Department of Biology, Concordia University, Montreal, Quebec H4B 1R6, Canada
- Supporting Information

ABSTRACT: Ocean metaproteomics is an emerging field enabling discoveries about marine microbial communities and their impact on global biogeochemical processes. Recent ocean metaproteomic studies have provided insight into microbial nutrient transport, colimitation of carbon fixation, the metabolism of microbial biofilms, and dynamics of carbon flux in marine ecosystems. Future methodological developments could provide new capabilities such as characterizing long-term ecosystem changes, biogeochemical reaction rates, and in situ stoichiometries. Yet challenges remain for ocean metaproteomics due to the great biological diversity that produces highly complex mass spectra, as well as the difficulty in obtaining and working with environmental samples. This review summarizes the progress and challenges facing ocean metaproteomic scientists and proposes best practices for data sharing of ocean metaproteomic data sets, including the data types and metadata needed to enable intercomparisons of protein distributions and annotations that could foster global ocean metaproteomic capabilities.

CHALLENGES IN METAPROTEOMICS

KEYWORDS: Metaproteomics, ocean, biogeochemistry, data sharing, best practices

INTRODUCTION

38 The measurement of many proteins within environmental 39 microbial communities, known as metaproteomics, is of 40 increasing interest to oceanographers and protein scientists. 41 The capacity to directly examine a multitude of functional 42 attributes of microbial communities and their linkages to both 43 ecology and biogeochemistry was once aspirational, but now 44 appears achievable with recent improvements in genomic 45 sequencing and mass spectrometry technology. Emerging 46 metaproteomic methodologies, in concert with other tradi-

tional and new approaches, could be particularly powerful in 47 the study of how complex environmental systems operate, as 48 well as how they respond to environmental changes.

Since the development of mass spectrometry based 50 proteomic technologies, there has been an increasing number 51 of metaproteomic or community-based analyses (Table S1), 52 including those of marine/ocean biota (Table 1). Metapro- 53 t1

Received: September 26, 2018 Published: January 31, 2019

Table 1. Examples of Ocean Metaproteomic Studies

North Atlantic Ocean, Bermuda Atlantic Time Series Station	Sowell et al., 2008; Bridoux et al., 2015; Saito et al., 2017
Ocean scale metaproteomics in the Atlantic Ocean	Morris et al., 2010; Bergauer et al., 2018
Antarctic Peninsula, Southern Ocean	Williams et al., 2012
Bering Sea Algae	Moore et al., 2012, 2014
Targeted metaproteomics of Central Pacific Ocean	Saito et al., 2014; Saito et al., 2015
Marine biofilms, shiphull environments	Leary et al., 2014
Metaproteomics of the Saniitch Inlet Oxygen Minimum Zone, Coastal Pacific Ocean	Hawley et al., 2014
Metaproteomics of aquatic estuary microbial communities	Colatriano et al., 2015
Marine sediments	Moore et al., 2012, 2012, 2014
Phaeocystis and diatom blooms in the Ross Sea of Antarctica	Bertrand et al., 2013; Bender et al., 2018

54 teomics of complex environmental samples such as seawater, 55 sediments, sinking particles, and biofilms have great potential 56 for revealing insight into biogeochemical cycling and microbial 57 response to environmental change in marine systems. For 58 example, recent ocean metaproteomic studies have provided 59 new insights into microbial nutrient transport, 1,2 colimitation 60 of carbon fixation processes, biogeochemical processes within 61 oxygen minimum zones, the composition of microbial 62 biofilms, dynamics of carbon flux in marine ecosystems, 6-8 63 and seasonal shifts in microbial metabolic diversity. Future 64 methodological developments should lead to new capabilities 65 such as characterizing large scale ecosystem changes,

estimating biogeochemical reaction rates from enzyme 66 concentrations and conducting in situ stoichiometric measure- 67 ments. In the short time since the emergence of these 68 metaproteomic methods, they have been applied to environ- 69 ments around the world: including coastal and open ocean 70 pelagic environments from the Atlantic and Pacific Oceans, 71 and even to the rapidly changing polar environments of the 72 Arctic and Antarctic regions. Diverse biological communities 73 have been sampled including free-living microbial and algal 74 communities (including microbiomes), sinking particles, 75 marine sediments, and even biofilms attached to human built 76 environments. 1-3,5,7,8,10-15 Also critical to the development 77 and deployment of metaproteomic approaches in natural 78 environments are controlled laboratory experiments on 79 cultivated microbes from the environment, 10,11,16-23 which 80 can enable the identification and verification of protein 81 biomarkers that characterize environmental processes.

CONFRONTING CHALLENGES IN METAPROTEOMIC RESEARCH

Despite this progress, key challenges remain in the application 85 of proteomic methods in environmental contexts.²⁴ These 86 challenges can be organized into four broad categories: (1) 87 environmental sample acquisition and protein extraction, (2) 88 chromatographic separation and mass spectrometry analysis, 89 (3) informatic data processing, and (4) data archiving and 90 sharing (Figure 1). A defining feature that affects all of these 91 f1 categories is that the ocean and other natural environments 92 contain a multitude of organisms that are not easily separated, 93 and hence are typically studied together in this "meta" 94

83

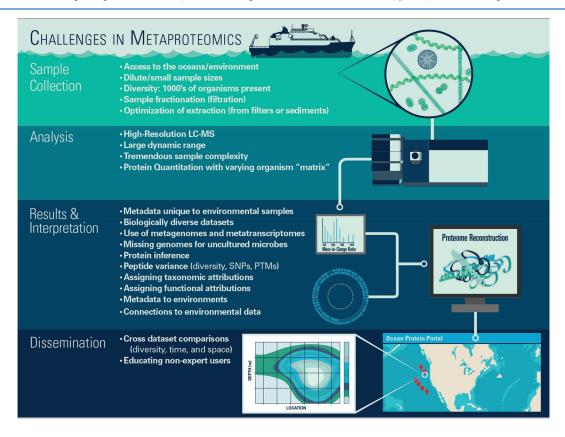


Figure 1. Analysis of proteins within natural environments presents unique challenges that can be improved upon to allow this new type data to inform ecosystem function and change. These challenges span sample collection and extraction, mass spectrometry analysis, informatic approaches, and data management and dissemination.

(a)



(b)

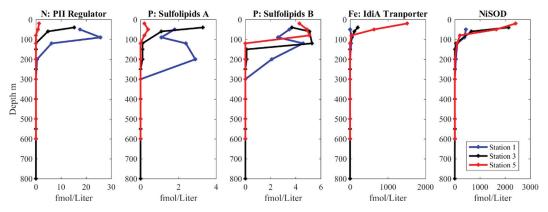


Figure 2. (a) Collection of ocean metaproteomic samples by in situ underwater McLane pump sampler as deployed in Terra Nova Bay of the Ross Sea in Antarctica aboard the icebreaker R/V Palmer to capture the microbial and algal communities as well as larger sinking particles by filtration of several hundreds of liters. (b) Example vertical distributions of three microbial proteins in the Equatorial Pacific Ocean using targeted metaproteomics that are biomarkers of nitrogen (N), phosphorus (P), iron (Fe) nutrient stress, and nickel (Ni) biogeochemical cycling (data from Saito et al., 2014, https://www.bco-dmo.org/dataset/646115). Proteins shown include the nitrogen PII regulator protein from Prochlorococcus (sequence VNSVIDAIAEAAK), the sulfolipid biosynthesis protein from Prochlorococcus (NEAVENDLIVDNK), UDP sulfolipid biosynthesis protein from multiple taxa (FDYDGDYGTVLNR), the IdiA iron transporter from Prochlorococcus (SPYNQSLVANQIVNK), and the nickel superoxide dismutase enzyme from Prochlorococcus and Synechococcus (VAAEAVLSMTK). Taxonomic assignments determined using METATRYP. 14

95 community context. For example, in a typical ocean seawater 96 sample, the microbial diversity includes prominent taxa from 97 each of the three major domains of life as well as from viruses. 98 This natural biological diversity manifests itself in a 99 tremendous chemical complexity for proteomics analysis, 100 where proteins from many organisms are digested into peptides and analyzed together, resulting in peptides that 102 have the potential to be shared across multiple species or ecotypes, or whose sequences are not within available DNA databases. New generations of fast scanning high resolution mass spectrometry instrumentation, such as orbitrap and timeof-flight instruments, now allow deep interrogation of these complex samples and the many low abundance or chimeric peaks within them, thereby improving and elevating the 109 confidence of identification. However, shared chemical similarities across this biologically diverse environment creates a number of challenges throughout the metaproteomic workflow. In this document we identify and describe the status of these challenges in order to enable researchers from environmental fields and beyond to focus efforts on resolving 115 them. In addition, we propose a set of best practices for current 116 and future data sharing for ocean metaproteomic data sets in 117 order for researchers to make maximal use of current and 118 incoming data sets. This effort is necessary to enable

interoperability and accessibility as this exciting new data 119 type becomes more widely adopted and to allow critical 120 temporal comparisons as the field evolves.

Sample Acquisition from Natural Environments and Protein Extraction

The study of natural marine communities presents significant 124 challenges in sample collection far beyond that involved in 125 laboratory-based studies. First, accessing the vast oceans that 126 cover 70% of the Earth's surface can require expeditions on 127 research vessels to reach remote oceanic locations. Second, in 128 seawater environments microbes are often 3-4 orders of 129 magnitude more dilute than model organism laboratory 130 cultures. For example, marine microbial populations can 131 range from 1000 to 100 000 cells per milliliter in seawater 132 compared to model microorganism cultures, such as 133 Escherichia coli that exceed a billion cells per milliliter. This 134 dilute cellular abundance in freshwater and marine environments requires filtration of tens to hundreds of liters of 136 seawater by combining multiple sampling bottles or using 137 specialized in situ underwater pumping systems to yield useful 138 quantities of protein for mass spectrometry analyses (Figure 139 f2 2a). Similarly, collection of sinking particles and sedimentary 140 f2 samples can require specialized sediment traps and coring 141

123

142 devices. There is considerable room for improvement in 143 engineering of sample collection as well as methodological 144 verification of sample handling processes, due to the combined 145 challenges posed by large geographical and depth space to be 146 sampled and the need to concentrate dilute biological material 147 without altering the proteomic signal within those samples. 148 Preservation of proteins at ambient temperatures appears to be 149 possible for some marine microbes using high salt RNA 150 preservatives, allowing in situ environmental samplers to be 151 designed and built, and time series to be taken. For example, a 152 commercially available RNA preservative was shown to 153 preserve proteins within cyanobacteria biomass at room 154 temperature for a month with no reduction in the number of 155 protein identifications, although periplasmic and extracellular 156 protein alkaline phosphatase was observed to be variable, implying loss during filtration.²⁵ This study of dissolved proteins and their role in biogeochemical cycling is also of 159 interest but will likely require separate sampling procedures to 160 concentrate them from seawater. There are new robotic 161 autonomous underwater vehicles (AUVs) being developed that 162 are specifically designed for proteomic sampling in natural 163 environments. For example, the Clio AUV incorporates recent developments in in situ pumping systems²⁶ to collect a suite of 165 discrete protein and other biogeochemical samples by vertically 166 moving and holding position at 16 depths over 6 km of a 167 vertical ocean water column. Integrated over typical ocean 168 expeditions, improvements in sampling efficiency allowed by 169 AUVs such as Clio will enable greatly increased sampling 170 depth resolution and geographic coverage of the vast ocean 171 basins.²⁷

When laboratory and environmental scientists interact, 173 confusion can arise from differing definitions/expectations of 174 biological replication. The scientific approach and objectives of 175 environmental sampling are distinct from laboratory experi-176 ments. There are clear differences between laboratory 177 experiments that can be easily replicated and sampling the 178 constantly changing natural environment. The challenge in 179 sample acquisition in marine metaproteomics described above 180 can preclude the collection of replicates; for example, 181 commonly used in situ pumps are tethered to a single wire 182 and deployed at predetermined depths and take several hours 183 to filter large volumes. Since the ocean is a fluid environment, a 184 second sampling deployment would collect a slightly different 185 water mass in space or time, depending on if the sampler was 186 placed adjacent on the vertical hydrowire or as a successive 187 sampling deployment after completion of the first sampling. As 188 a result, real variations (albeit small) in biological communities 189 and chemical properties could be captured in attempts at 190 sampling replication, and true biological duplicates are aspirational, if not impossible. In place of replication, 192 oceanographers often look for "oceanographic consistency" 193 in trends across vertical depth structure (or horizontal 194 structure in the case of ocean basin sections) as a useful 195 means to validate results.²⁸ Single samples have demonstrated 196 this oceanographic consistency in capturing large scale oceanographic and metabolic processes across chemical and biological gradients.^{2,3,29}

The comprehensive extraction of proteinaceous material from biomass is another challenge in metaproteomic studies. Environmental samples can be extraordinarily complex due to being composites of significant biological diversity, as well as having additional biogenic and nonbiogenic materials within them. Moreover, the biological composition of metaproteomic

samples can be largely unknown prior to extraction. Hence, the 205 ability to tailor and optimize extraction protocols to the 206 environmental sample type presents unique difficulties. In 207 water column environments, depending on the environment 208 and collection strategy, an environmental microbial sample will 209 contain dozens of major biological species and hundreds to 210 thousands of trace species. 30,31 Sediment and sinking particle 211 samples contain not only a mixture of organisms, but partially 212 degraded peptides created by a phalanx of microbial proteases 213 produced by heterotrophic bacteria consuming those particles. 214 There are also numerous complex symbiotic communities such 215 as corals, hydrothermal vent tube worms, and other symbiotic 216 systems where the proteins of the microbial assemblage will be 217 present within the extensive proteome of a eukaryotic host. 218 Studies have examined the recovery efficiency of different 219 extraction buffers on sedimentary and microbial biomass. 25,32 220 Moreover, the presence of biogenic soft and hard parts, 221 including mucilage, calcium carbonate, and siliceous compo- 222 nents, as well as mineral phases, can complicate chemical 223 separation of proteins and impair protein extraction efficiencies 224 and require development of matrix-specific extraction proto- 225

Mass Spectrometry Analyses

To date, the mass spectrometry measurement component of 228 metaproteomics has utilized three types of approaches: data-229 dependent acquisition (DDA) for discovery proteomics, 10,34,35 230 data-independent acquisition (DIA^{2,36}), and targeted meta-231 proteomics for quantitative analysis using multiple or parallel 232 reaction monitoring approaches (MRM/PRM^{3,11,14}). 233

Briefly, these approaches differ in how they select ions for 234 fragmentation: DDA approaches continually select abundant 235 features within ms¹ spectra for further ms² fragmentation 236 analysis (isolating the most abundant peaks within each parent 237 ms¹ spectra for fragmentation, with various user parameters 238 such as excluding recently fragmented precursors for a short 239 period),³⁷ while DIA methods conduct ms² fragmentation on 240 small sequential mass windows across the entire mass range of 241 interest,³⁸ thereby potentially fragmenting spectra of all ions, 242 assuming sufficient intensity. In contrast, targeted methods 243 focus their fragmentation analyses only on precursor ions 244 found on the target list, thereby increasing sensitivity by 245 focusing mass spectrometry time on target ions. 39-41 DDA 246 approaches continue to be most prevalent in metaproteomics, 247 but targeted and DIA approaches are increasingly being 248 explored for their ability to provide absolute and relative 249 quantitation, respectively. An example DDA workflow is shown 250 in Figure 3, and examples of vertical profiles of targeted 251 f3 peptides from MRM/PRM experiments are shown in Figure 252 2b.

While these proteomic methods have become common in 254 proteomic analyses on single organisms, the complexity of 255 metaproteome samples presents challenges for each method 256 with regard to both the chromatographic separation and mass 257 spectrometry components. For comparison, the complexity of 258 ocean seawater metaproteome samples appears to be 259 significantly greater than the human proteome, despite the 260 latter typically being considered to be one of the more complex 261 proteome sample types. This is illustrated in Figure 4a where a 262 f4 three-dimensional (3D) visualization of the mass spectra 263 acquired from a surface sample in the central Pacific Ocean is 264 shown (filtered by 0.2–3.0 μ m size fraction range), and in 265 Figure 4b–c with spectra from a small mass range examined at 266

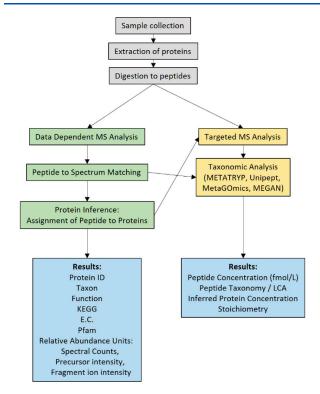


Figure 3. An example environmental metaproteomic workflow where environmental samples are collected and extracted (gray), discovery proteomics are conducted (green), and peptide targets from selected proteins of interest can be assayed using isotopically labeled peptide standards whose taxonomic assignment can be queried against databases of genomes and metagenomes (yellow). The results can provide relative and absolute abundance measurements of the protein from the microbial and algal community, including functional and taxonomic information (blue).

267 equivalent chromatographic elution times $(575-578 \text{ m/z ms}^{1})$ 268 window and 140-141 min) revealing more observable mass 269 peaks events in an ocean sample (Figure 4c) when compared 270 to a human cell line (HeLa) sample (Figure 4b). These 271 observations of metaproteome complexity were also quantita-272 tively confirmed across entire samples by analysis of ms¹ peak 273 within triplicate HeLa injections and five metaproteome 274 samples from the Pacific Ocean at varying depths in Figure 275 5. These HeLa-ocean comparisons used identical chromato-276 graphic and mass spectrometry settings and were run within 1 277 week of each other using the same nanospray column, with 0.5 $_{278}$ μg of HeLa analyzed per injection, while 1.0 μg ocean sample 279 was analyzed per sample injection. In this example, the number 280 of peaks was higher in the metaproteomes compared to HeLa 281 (Figure 5a-c), while the total ion current (TIC) was 282 considerably lower across all metaproteome samples (Figure 283 5a-b,d), implying more peaks of lower intensity in the 284 metaproteome samples. This high complexity of metaproteo-285 mics samples presents significant challenges to current 286 chromatographic and mass spectrometry workflows. For 287 example, real-time feature identification (peak picking) 288 software on mass spectrometers has not been optimized to 289 process chimeric peak features that appear to be ubiquitous in 290 metaproteomic samples, where chimeric features are peaks so 291 close in mass to other peaks preventing a successful charge 292 state estimate that is needed to trigger ms² fragmentation in 293 bottom up DDA experiments. Moreover, the low abundance of 294 many ions in metaproteomic samples (as observed in Figure

4c) poses an additional challenge, where the numerous low 295 abundance peaks among more abundant ones remain 296 uncharacterized due to physical limits on the number of ions 297 entering the mass spectrometer at any time, a problem that can 298 challenge both DDA and DIA methods.

Metaproteomic approaches have made progress in address- 300 ing the challenges of this sample complexity. For example, 301 chromatographic approaches have been improved by applying 302 two-dimensional chromatography^{10,34} or gas phase fractiona- 303 tion^{2,7,36,43} to distribute the sample complexity for mass 304 spectrometry analysis across subsamples or temporal chroma- 305 tographic separation as a means to obtain deeper metapro- 306 teomes. Moreover, DIA approaches have been utilized to 307 address the crowded and complex nature of ion chromato- 308 grams that are specific to metaproteomics,^{2,36} although 309 bioinformatic pipelines for mixed community DIA data sets 310 are still being developed. Finally, the application of targeted 311 methods offers improved sensitivity and absolute quantitation 312 of biomarkers for environmental stress by targeting represen- 313 tative peptides.^{3,11,14}

Future collaboration with hardware and software developers 315 could also greatly improve metaproteomic research efforts. For 316 example, effort could be expended to capture greater 317 information about the numerous low intensity ions that are 318 missed by real-time and postprocessing algorithms due to 319 several factors including insufficient ions trapped for high- 320 quality ms² fragmentation spectra, ions being chimeric with 321 other nearby peaks, and lack of charge state assignments. 322 Recent efforts in improving detection of chimeric peaks may be 323 useful in this regard when applied to metaproteomic 324 applications. 44

Finally, there is an important need for intercomparison and 326 intercalibration efforts with regard to protein extraction 327 efficiency and mass spectrometry accuracy and precision. 328 Chemical oceanographers have a legacy of successful 329 intercalibration efforts that have enabled global scale studies 330 of ocean chemistry, such as the recent GEOTRACES (an 331 international study of the marine biogeochemical cycles of 332 trace elements and their isotopes) trace elements and isotope 333 global section program.⁴⁵ For ocean metaproteomics, uniform 334 preparation of large batches of intercalibration samples may be 335 challenging given that samples can vary in biological 336 composition and sampling methodologies, and likely multiple 337 smaller initial intercomparison studies might first be needed. 338 Alternatively, simpler "synthetic" metaproteome samples could 339 be created by mixing of laboratory microbial isolates that could 340 be made in large batches and distributed, although these may 341 not reproduce the depth of biological diversity nor a realistic 342 environmental chemical matrix. Intercalibrations could be 343 applied to the two current major approaches to ocean 344 metaproteomic mass spectrometry analysis: global discovery 345 data sets and targeted metaproteomics, with studies providing 346 metagenomic databases and isotopically labeled peptide 347 standard materials to facilitate analyses, respectively. Moreover, 348 intercalibration exercises could be conducted on consensus 349 standard sample sets of some example biological communities 350 initially, such as seawater microbial communities that are well- 351 characterized with respect to metagenomic data, although 352 eventually many types of biological materials could be selected 353 for intercalibration (sediments, biofilms, etc.). Finally, future 354 additional types of metaproteomic analyses could be added for 355 intercomparison such as data independent analysis and post- 356

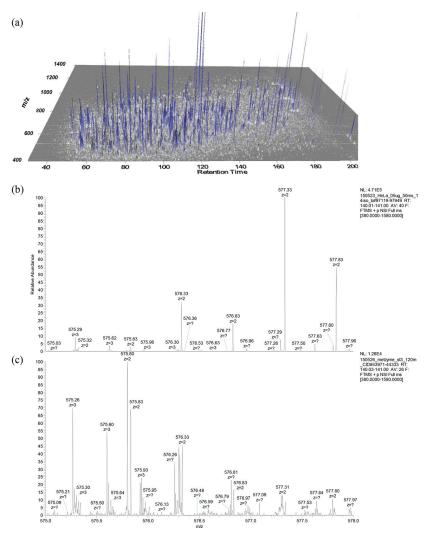


Figure 4. (a) Three-dimensional representation (axes of retention time (min), m/z, and intensity) of complex spectra associated with an environmental ocean sample from the Equatorial Pacific from the METZYME expedition (200 m depth, produced in MzMine2⁴²). Comparison of a 3 m/z ms¹ mass window (575–578 m/z, 140–141 min) from (b) human proteome spectra (HeLa cell line) and (c) ocean metaproteome (120 m depth) provides an example of the high complexity of environmental samples due to the biological diversity.

357 translational modifications within metaproteomic environ-358 mental samples.

359 Metaproteomic Data Analysis

360 Data analysis of mass spectra from metaproteomics experi-361 ments presents many challenges compared to single organism 362 proteomics. In particular, each metaproteome mass spectral 363 data set can contain tremendous biological diversity whose 364 composition is often largely unknown. Furthermore, established proteomic workflows that conduct peptide-to-spectrum matching (PSM) by comparing peptide precursor and fragment ion masses to corresponding predicted masses using genomic reference databases were never designed to handle the inherent complexity and multiple biological entities 370 within metaproteomic data sets, and hence approaches thus far have been improvised adaptations. The expanse of unknown 372 biological diversity often results in metaproteomic protein 373 database searches that are typically large and of redundant 374 nature. This has an effect on database selection, data-search 375 algorithm utilized, subsequent FDR statistics, 46,47 and protein 376 inference. 48,49 Additionally, metaproteomics shares the chal-377 lenge of functional and taxonomic assignments with metagenomics, relying on a comparative approach with model 378 organisms, resulting in many proteins of unknown function or 379 taxon. Finally, metaproteomic workflows typically involve the 380 integration of multiple software tools, making documentation 381 and reproducibility difficult as tools evolve. Despite these 382 challenges, multiple approaches have been developed over the 383 last 13 years (Table S1). The analytical workflows that have 384 been developed to date are mainly comprised of (a) database 385 generation, (b) database search, and (c) taxonomy and 386 functional analysis.

Database Type (Genome, Metagenome, Metatran- 388 scriptome, Custom). In order for PSM algorithms to assign 389 peptide sequences to spectra from MS experiments, the 390 observed tandem mass spectra are cross-correlated and scored 391 against theoretical spectra generated in silico from the 392 provided protein sequences. The collection of protein 393 sequences is generated from available genomic, metagenomic, 394 or metatranscriptomic sequence information and commonly 395 referred to as the genomic or protein database. High scoring 396 peptide spectral matches (PSMs) are then reported with their 397 corresponding protein sequence and annotation from the 398 original database. Importantly for metaproteomics, the 399

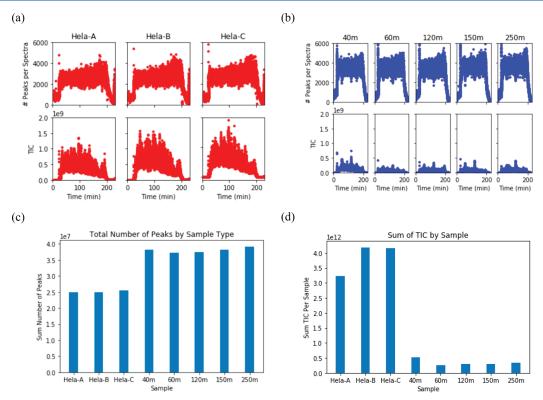


Figure 5. Peak analysis of human cell line and ocean metaproteome samples by identical chromatographic and mass spectrometry conditions. (a–b) Number of peaks identified replicates (top rows) and the total ion current (TIC, bottom rows) of the sample in Hela (Panel A, replicates Hela-A, Hela-B, and Hela-C) and an ocean metaproteome sample (Panel B, depth 40, 60, 120, 150, and 250 m, Metzyme Expedition Station 3). Samples were run during the same week on the same nanospray column (see methods) with similar amounts of protein injected (0.5 μ g for Hela per injection, 1 μ g for ocean metaproteome). (c) Total number of peaks by sample type showed a higher number peaks in ocean metaproteome samples, while (d) TIC by sample showed much lower summed peak intensity within the metaproteome samples, consistent with the 575–578 m/z window shown in Figure 4.

400 peptides and proteins reported are dependent on the 401 coherence of the original genomic sequence information relative to the organism(s) present in the sample. More often than not in metaproteomics, each sample's composition of biological diversity is unknown or its characterization is 405 limited by the depth of DNA sequencing and final assembly. As 406 a result, if a peptide sequence is not in the database, the 407 peptide will not be identified nor will its contribution to a 408 protein identification be included in the experiment. 409 Furthermore, quality of gene prediction algorithms can affect 410 protein detection: if protein-encoding genes are missed during 411 the initial gene prediction phase, then they will not be included in the protein search database. While gene prediction from prokaryotic genomes is relatively straightforward, it becomes 414 challenging for more complex microbial eukaryotic genomes, 415 owing to the complexity and diversity of eukaryotic gene 416 structure (e.g., predicting introns and exons). However, eukaryotic gene prediction algorithms are continually advancing, and indeed proteomics plays a large role in the accurate 419 identification of protein encoding regions of eukaryotic genomes through proteogenomic efforts. 50-52 Additionally, the incomplete nature of peptide fragmentation yields high variability in final peptide interpretations, making database choice and construction pivotal.⁵³ Finally, the occurrence of 424 similar but not identical protein sequences (homologues) in 425 closely related organisms adds significant complexity to 426 metaproteomic search databases.

There are three main approaches for creating metaproteomic 427 databases: (1) sequence and assemble a metagenome, (2) 428 assemble a database from the public environmental genomic 429 repositories, and (3) create a pseudo-metagenome by including 430 desired taxonomic classes or species. The composition of the 431 protein search database used to search the mass spectra from a 432 metaproteomic sample has a profound effect on biological 433 conclusions.⁵⁴ Timmins-Schiffman et al. recommended a best 434 practice for environmental proteomics of corresponding site 435 and time specific metagenomes in order to generate accurate 436 proteomic databases to assign peptide sequences and protein 437 annotations. 3,4,10,55,56 While this avenue represents the ideal 438 scenario, at some point sufficient metagenomic coverage of 439 specific environments should allow decoupling between 440 genomic and proteomic analyses as a large inventory accrues 441 of deeply sequenced data sets from diverse environments.⁵⁷ 442 However, as evolution is a dynamic process, resequencing of 443 these environments will be required to capture continued 444 community adaptation to changing environments and evolu- 445 tionary forces which are already evident in repeated marine 446 sequencing efforts over seasonal time scales.⁵⁸

There are a variety of publicly available metagenomics data 448 sets that marine metaproteomics researchers have used, for 449 example, the J.C. Venter Institute's Global Ocean Sampling 450 (GOS) database. ^{1,59-61} In addition, there are environmental 451 metagenomics databases available at major repositories and 452 portals such as EBI, JGI, and iMicrobe (https://www.ebi.ac. 453 uk/metagenomics, https://img.jgi.doe.gov, http://www. 454

455 imicrobe.us). For eukaryotic phytoplankton and protists, 456 genomic, transcriptomic, and metagenomic resources are 457 considerably more scarce, though recent availability of the 458 Marine Microbial Eukaryote Transcriptome Sequencing 459 Project (MMETSP) has begun to address this challenge. 62 460 The application of large databases (either public or 461 metagenomics) still suffers from limitations with respect to 462 sensitivity of identifications. One approach to alleviating this 463 problem applied a "metapeptide database" from shotgun 464 metagenomics sequencing and demonstrated a significant 465 increase in the number of identifications presumably due to 466 a more accurate and compact database as compared to an 467 assembled predicted metaproteome and NCBInr. 63

Finally, the selection and compilation use of individual 469 microbial genomes for a metaproteomic database can also be 470 useful in metaproteome analysis. Given that many of the major 471 microbial taxa in the oceans were discovered in recent decades 472 and in many cases there are few laboratory isolates and 473 accompanying genomes, there is significant need to amend 474 large public databases with new representative microbial 475 genomes. 64-66 In contrast to metagenomic data sets, these 476 genomes of cultivated isolates also provide clarity with regard 477 to taxonomic attribution that can be obfuscated by limitations 478 of metagenomics assembly and annotation. Increasing 479 availability of single cell genomic data (single amplified 480 genome; SAGs) can also contribute significantly to databases 481 for metaproteomic analysis. Notably however, the SAG 482 technology does not produce complete genomic sequences 483 (unclosed genomes), and hence care must be taken when 484 trying to interpret the absence of a protein using metagenomics 485 or SAG databases to avoid false negatives. For eukaryotic 486 metaproteome analysis, transcriptome data can also serve as a 487 useful source of sequence for the protein database generation 488 since full genomic information of marine eukaryotes is 489 relatively rare and the DNA architecture is more complex 490 due to the prevalence of noncoding intron regions intertwined 491 with the protein-coding exons. A recent study from the Ross 492 Sea of Antarctica found that for a diverse bloom community 493 with abundant eukaryotic phytoplankton, a combined database 494 of transcriptomes from cultured isolates and field metatran-495 scriptome provided a richer metaproteome result than either 496 database alone. 10

Search Engine. In common DDA proteomic workflows, 498 the search engine that conducts the PSM analysis is central to 499 protein discovery and identification. Application of these PSM $_{500}$ algorithms (e.g., SEQUEST, X!Tandem) $^{67-69}$ have been 501 successfully applied to metaproteomic analyses, despite the 502 fact that they were never designed to deal with the complexity 503 of metaproteomic data sets. Search algorithms are chosen 504 based on the following factors such as the ability to search large 505 databases, speed, and the ability to generate outputs that are 506 compatible with downstream processing steps such as peptide 507 or PSM output with robust FDR threshold calculations. Most 508 of the suggested database generation strategies generate large 509 databases, which in turn affect the sensitivity of identifications. 510 Multiple strategies have been suggested to increase peptide identifications. This includes the two-step method for searching large databases; 70-74 and a cascaded search method. Muth et al. have recommended using a database 514 sectioning approach so that searches against subsets of a large 515 database may increase the number of identifications. 53 They 516 have also suggested using multiple search algorithms in order 517 to increase the percentage of peptide spectral matches in a data

set. For example, SearchGUI/PeptideShaker,⁷⁶ which uses at 518 least eight open-source search algorithms, can facilitate this 519 multipronged approach and can be used to search against large 520 databases.⁵³ Irrespective of the choice of search algorithms, the 521 goal is to generate outputs with maximal coverage of mass 522 spectra that are compatible with the next steps of taxonomic 523 analysis, functional analysis, and subsequent targeted valida- 524 tion

Despite these initial successes, it is apparent that these 526 workflows and algorithms could be improved upon to confront 527 significant challenges of spectral complexity, metaproteomic 528 protein inference, and taxonomic attribution within environ- 529 mental samples. Specifically, the presence of numerous low 530 abundance peak features, discerning chimeric peaks (Figure 2), 531 and assignment of corresponding peptide charge states are 532 difficult for current PSM algorithms and likely result in 533 significant underestimation of peptide identifications within 534 metaproteomic spectra. Application of de novo search 535 algorithms and spectral libraries could also improve identi- 536 fication of peptides from metaproteomics samples. 537

Taxonomic and Functional Annotation. Metaproteo- 538 mics has a distinct utility in determining the protein functional 539 expression by a microbial community.⁷⁴ However, the 540 functional interpretation of a metaproteomics data set is 541 inherently reliant upon the underlying annotation of the 542 protein search database, including the prediction of protein- 543 encoding genes from genomic data and the subsequent 544 functional annotation of the predicted proteins. While much 545 of the taxonomic and functional attribution of metaproteomic 546 results can leverage metagenomic annotation pipelines, there 547 are aspects that are unique to metaproteomics. In particular, 548 the basal unit of proteomic identification is generally the 549 tryptic peptide (due to the effectiveness of trypsin in 550 proteolytic digestions), resulting in amino acid sequence 551 coverage without overlaps, except in cases of missed cleavages. 552 Due to the presence of unknown biological diversity, it is 553 possible to have tryptic peptides that are shared within or 554 across species. As a result, the greatest confidence in 555 metaproteomic discovery occurs on the peptide level, and 556 creates a need in metaproteomic research for investigation of 557 sequence taxonomy on the peptide level. This is also an issue 558 since inference of specific taxonomy tends to be more difficult 559 than function in typical sequence analysis (e.g., BLAST), due 560 to the sharing of biochemical capabilities by many organisms. 561 Two web-based applications are available for this, Unipept and 562 Metatryp, that search DNA sequence data for the presence of 563 user entered tryptic peptide sequences and estimate the lowest 564 common ancestor (LCA; kingdom, phylum, genus or species) 565 for each peptide query (Figure 2b). The applications are 566 distinct in the DNA sequence databases they search, with 567 Unipept searching the UniProt genomic database as well as 568 providing cross-referenced EC numbers and GO terms, 80 while 569 Metatryp allows use of custom genomic data and metagenomic 570 data (including single amplified genomes and metagenome 571 assembled genomes) with a focus on marine environments 572 (http://metatryp.whoi.edu).14 The choice of database can 573 affect results; for example, currently Unipept maps 51% of the 574 peptides from the Morris et al. South Atlantic data set² to 575 sequences within Uniprot, implying that genomic data 576 availability still hinders interpretation of ocean metaproteomic 577 data sets (https://unipept.ugent.be/mpa). There are addi- 578 tional bioinformatic tools for taxonomic analyses that may be 579 useful for metaproteomic research such as MEGAN⁷⁶ 580

581 microbiome software that computes taxonomic profile by 582 assigning PSMs from a metaproteomics experiment to an 583 appropriate taxonomic unit within the NCBI taxonomy. In 584 addition, the recently developed MetaProteomeAnalyzer that 585 uses outputs from SearchGUI/PeptideShaker⁸¹ has taxonomic 586 analysis capability.

Connecting protein functions with metaproteomic data sets 588 is a key goal that can be accomplished in variety of ways. 589 BLAST analyses of the metagenomics contigs being used for 590 PSMs provide high-quality searches by using longer sequences, 591 but require availability of well-annotated metagenomics 592 databases. Additional software is available for downstream 593 metaproteomic functional analysis including peptide-level MetaGomics⁸² or Unipept,⁸³ protein-level (for example, 595 MEGAN⁸⁴), protein orthologs (for example, EggNOG 596 mapper⁸⁵ or metaprotein/protein-group level (for example, 597 MetaProteomeAnalyzer). Each of these methods uses distinct 598 annotation databases, such as UniProt (for example software 599 tools such as MetaGomics, Unipept or MetaProteomeAna-600 lyzer⁵³) or NCBInr (MEGAN) or EggNOG database 601 (EggNOG mapper) to assign functional categories. Functional 602 analysis tools generate functional ontologies such as Gene 603 Ontology (GO; for example, MetaGOmics, MEGAN, Unipept 604 and EggNOG mapper), KEGG orthology groups (for example, 605 EggNOG mapper, MEGAN and MetaProteomeAnalyzer), EC 606 numbers (for example, Unipept and MetaProteomeAnalyzer), 607 and EggNOG orthologous groups (EggNOG mapper) are 608 used for deciphering the functional state of a microbiome.

While these annotation approaches described above are 610 useful, it is worth acknowledging that there is a continuing 611 challenge in interpretation of metaproteomics data that is 612 inherited from metagenomic research regarding the process of 613 annotating protein function from gene sequence data. The vast 614 majority of protein annotations are assigned not from direct 615 experimental evidence, but rather from sequence similarity to a 616 previously annotated protein or protein family in a metabolic/ 617 ortholog database (e.g., KEGG, COG, METACYC, PFAM). 618 This leads to several issues. The first is that annotation transfer 619 based upon sequence similarity has resulted in the propagation 620 of misannotations in large genomic databases over time, with 621 the most common form of misannotation resulting from "over 622 annotation"—annotation of a gene to a deeper level of 623 functional characterization than the supporting evidence 624 provides. 86 Even minor errors or discrepancies in annotation 625 transfer can result in massive error propagation.⁸⁷ Common 626 irregularities in gene annotation can cause serious issues for the 627 metaproteomics researcher who is reliant upon these 628 annotations to give biological context to proteomic data. 629 Some of these annotation irregularities include gene 630 annotations in which one gene name is assigned very different 631 functional descriptions, instances where the gene name for a 632 particular function has changed over time, or cases where only 633 one function of a multifunctional enzyme is provided. 87 634 Similarly, novel functions can be discovered for previously 635 unannotated hypothetical proteins. 16,88 These issues can be 636 compounded if a custom search database comprised of 637 genomes annotated by different means is used for peptide 638 and protein identification, as is common in oceanographic 639 studies. Moreover, as genomic data are updated with improved 640 annotation information, an ability to pass this new information 641 onto deposited processed metaproteomic results will be 642 needed, and could be accomplished with versioning of 643 deposited data sets.

A barrier to managing the spread of misannotations is that 644 for some databases, such as GenBank NR, there are currently 645 no means for the community to submit additional manually 646 curated annotations and corrections. Fortunately, newer 647 techniques for genome annotation which rely on methods 648 beyond simple pairwise sequence similarity—most notably, the 649 use of machine learning based algorithms—outperform the 650 former pairwise similarity and BLAST based annotation 651 transfer methods. 89 Protein functional prediction from 652 sequence data is a growing field itself, and will likely benefit 653 from coupling powerful predictive algorithms with high-quality 654 systems data to provide deeper, more accurate, and more 655 meaningful characterizations.

Another common issue is where only a single function is 657 reported for a protein family that is comprised of proteins of 658 divergent function. For example, in Colatriano et al., 56 659 numerous proteins are assigned to the DMSO reductase 660 enzyme superfamily, which is comprised of a number of 661 functionally distinct proteins including nitrate reductases 662 involved in anaerobic respiration as well as nitrite oxidor- 663 eductase involved in dissimilatory nitrite reduction. Only 664 through fine-scale phylogenetic analysis of the identified 665 proteins could the true function as nitrite reductases be 666 determined. However, in many cases, the relationship between 667 phylogeny and function within protein families is unknown. In 668 metaproteomics, this is especially problematic for transporter 669 proteins which are often abundant in metaproteomics data sets 670 and are attractive since they may be used to infer substrate 671 utilization patterns that are directly relevant to global 672 biogeochemical cycles. However, transport function is often 673 poorly conserved with these families, and hence sequence 674 analysis is no substitution for the critical biochemical and 675 genetic studies capable of characterizing protein function 676 directly. 16,90 In conclusion, there are significant challenges and 677 room for improvement in the assignment of annotation 678 information to metaproteomic data sets, and cooperation 679 with genomics researchers and organizations, as well as 680 incorporating an ability to reanalyze data sets and submit 681 updated versions will be important components of future 682 metaproteomic data management.

Challenges in Ocean Metaproteomic Data Sharing

There is potentially great value in sharing raw and processed 685 environmental metaproteomics data within ocean sciences and 686 beyond. As with most 'omics sciences, each data set contains 687 far more information than the data-generator's laboratory can 688 interpret. Proof that a gene is synthesized into protein form, as 689 well as its variation in spatial or temporal distribution, can 690 provide valuable biological and chemical information about the 691 environment. Yet due to the complexity and newness of this 692 data type, there are challenges unique to metaproteomics in 693 reporting and disseminating this information. In a workshop in 694 2017, we discussed these challenges, and organized the 695 following information in an attempt to provide a first set of 696 best practices to Ocean Metaproteomics Data Sharing.

Interoperability between ocean metaproteomic observations 698 and their related environmental data requires that the 699 relationships between these data are explicitly known. Defining 700 these relationships helps to communicate proper use when 701 integrating disparate data within a shared domain.⁹¹ While 702 defining these relationships helps humans properly integrate 703 data, software and tools need something more. In the past, this 704 meant developing software to a specific set of data types that 705

I

706 forced data to follow certain conventions for variables names 707 and structure. Yet current standards for integrating data on the 708 web enable software to infer how disparate data can be 709 integrated when they are described using semantic web 710 schemas and rulesets. In doing so, these disparate data do 711 not need to be transformed to conform to the software. 712 Instead, the data are described using semantic web 713 technologies for proper integration. The semantic web 714 provides a framework for classification of data and its 715 relationships in what are called ontologies, or vocabularies. 716 These ontologies are logical groupings of terms and the axioms 717 that define the how data from that domain are to be described. 718 These terms can define cardinality rules or other logical 719 expressions that enable humans and machines to make 720 inferences over the data. Instead of transforming the original 721 data to force it to conform to software, software can be written 722 to conform to an ontology. As a result, this leaves the data 723 intact and moves the work of integrating data to describing 724 how it maps to the ontology. 92

725 CONNECTIONS TO KEY ENVIRONMENTAL 726 METADATA AND DATA

727 Environmental research requires unique metadata to provide 728 context for comparisons across space and time. These 729 metadata include numerous attributes associated with sampling 730 from the oceans or other natural environments that are most 731 often not included in the data model in biomedically focused 732 proteomics repositories. For example, there is geospatial and 733 environmental contextual information that is critical to 734 interpreting results such as latitude and longitude, depth, and 735 sampling environment (e.g., pelagic water column or benthic 736 sedimentary location). For the pelagic environments, there is 737 critical methodological information regarding sample collec-738 tion parameters including filtration pore size range(s) or 739 sediment trap deployment conditions for sinking particles. For 740 benthic environments, key sedimentary environmental details, 741 organism (coral, whale, plankton, zooplankton, etc.), or human 742 built environments (ship hull surface) need to be described. In 743 addition, local time of sampling can be important to detect 744 short-term (diurnal) changes or long-term (seasonal or 745 environmental change) processes. In addition to these 746 metadata, there is also important contextual information 747 derived from co-occurring chemical, biological, or physical 748 measurements, such as temperature, macronutrient and 749 micronutrient abundances, salinity, light intensity, and bio-750 logical diversity, to name a few parameters. Carefully defining 751 the data and metadata model will also facilitate connections to 752 environmental data management holdings such as those at the 753 Biological and Chemical Data Management Office for Ocean 754 Science in the US (www.bco-dmo.org), and various national 755 data repositories will facilitate access to this contextual 756 information. Table 2 provides a list of recommended metadata for best practices of ocean metaproteomic samples data management to be provided at the time of deposition.

Making connections between metaomics data sets and roo environmental data is a widely sought goal that is difficult to achieve. Enabling interoperability between ocean metaproteomic observations and their related environmental data requires that the relationships between these data are explicitly known. Defining these relationships helps to communicate proper use when integrating disparate data within a shared domain. While defining these relationships helps humans properly integrate data, software and tools need something more. In the

Table 2. Recommended Metadata for Best Practices in Ocean Metaproteomics Data Sharing

1	8	
reporting metric	notes/units	when available
Project Metadata		
expedition identifier		
lead PI, contact info, ORCID identifier		
Co PI, contact info, ORCID identifier		
Contextual Metadata		
latitude	degrees N	
longitude	degrees W	
sampling depth		
habitat type	Pelagic, benthic, reef, ship-hull, host-associated, other	
temperature	degrees Celsius	
salinity		$\sqrt{}$
chlorophyll-a concentration		√ √ √
oxygen concentration		$\sqrt{}$
other analytes measured	links to environmental data repositories	
Sample Acquisition Metadata		
sampling method	filtration, sediment trap, coring, other	
volume of water sample represents	liters	$\sqrt{}$
filter type	membrane (PC, PS), glass fiber, quartz, other	$\sqrt{}$
filter size	micron pore size	\checkmark
prefilter(s) used	if applicable: micron pore size, filter type	$\sqrt{}$
Sample Extraction Methods	see Table S2	
Mass Spectrometry Methods	see Table S2	
Data Analysis Methods	see Table S2	
Metaproteomics Data Analysis Metadata		
database used for PSM or targeted method development	metagenomic, metatranscriptomic, genomic sequences used; link to sequence repository	
taxonomy analysis method	software/algorithm used	
functional analysis method	software/algorithm used	

past, this meant developing a certain piece of software to a 768 specific set of data types. Yet current technologies, such as the 769 semantic web, enable software to understand how data can be 770 integrated through well-defined schemas and rulesets. Using 771 ontologies, a semantic web technology, data, and their 772 necessary relationships can be described in ways that machines 773 can enforce cardinality constraints and make inferences that 774 are helpful for ensuring a proper integration. 92

Due to this fundamental importance of metadata associated 776 with metaproteomics results, deposition of raw data into 777 existing proteomic repositories designed primarily for labo- 778 ratory studies (e.g., Pride, Massive, and Chorus) could create 779 challenges for researchers in locating and manipulating 780 collections of environmental data sets. While these raw 781 data repositories are already valuable in hosting environmental 782 proteomic data, the proper data management of large amounts 783 of metadata can be viewed as a burden beyond what those 784 entities are funded to provide, as has been observed in 785 metagenomics data management spheres. As a result, 786 intermediary environmental metaproteomics portals that host 787

Table 3. Key Datatypes Needed for Ocean Metaproteomic Data Sharing

data type	attributes
raw mass spectra files	open data format, parameters, corresponding environmental and mass spectrometry metadata (see Tables 2 and S2)
protein identifications	MS/MS sample name
	sequence identifier (e.g., metagenome locus or genome ORF)
	product name, taxon, taxon ID, KEGG, E.C., PFam
	quantitative value(s) (spectral count)
peptide identifications	MS/MS sample name, sequence identifier (e.g., metagenome locus or genome ORF), peptide sequence, peptide start index, peptide stop index, precursor mass, retention time, statistical score of peptide, quantitative value (spectral count, MS¹ peak area, fragment ion peak areas from SRM/MRM/PRM analyses, calibrated SI unit concentrations)
FASTA of amino acid sequences of all identified proteins	full sequence of DNA or RNA used for peptide-to-spectrum mapping
	corresponding sequence identifier (e.g., metagenome locus or genome ORF)

788 processed data sets and full metadata archives can serve a 789 valuable function as a link to raw data repositories and 790 cocollected or colocated environmental data sets. Hopefully 791 either raw data repositories will work with environmental 792 communities in collecting environmental metadata as well as 793 cooperating to enable web-based connections between raw 794 data repositories and processed data portals. In order to foster 795 a high level of data sharing, reanalysis, and intercomparison, it 796 is important for data generators to preserve a number of data 797 facets and metadata, at the time of publication and archiving.

DATA TYPES USEFUL FOR METAPROTEOMIC DATA SHARING

800 In addition to the unique and critical metadata and 801 environmental data that need to be provided, the metapro-802 teomic data sets also require several distinct types of raw and 803 processed data (see Table 3) in order to allow reproducibility 804 and to enable a deep interrogation of their attributes in 805 environmental data portals such as the future Ocean Protein 806 Portal. Within processed data sets these include protein 807 identifications with associated functional and taxonomic 808 information (if known), full amino acid sequences, correspond-809 ing peptide sequences of discovered peptides, quantitative 810 information for both proteins and peptides (e.g., spectral 811 counts, precursor or fragment intensities), and associated 812 statistical threshold used for generating these data (e.g., protein 813 and peptide FDRs). For raw data, these include the raw mass 814 spectrometry files converted to a platform-independent format 815 as well as the sequence databases used to generate them. An 816 important distinction from model organism studies is the 817 fundamental importance of the peptide-level data to 818 metaproteomics: the ability to have access to detailed 819 peptide-level information with corresponding geospatial and 820 temporal information will be critical to enabling users to 821 directly interrogate the peptides that were actually measured in 822 the oceans, as opposed to relying on protein inference that may 823 be incorrect due to insufficient metagenomic coverage.

QUANTIFICATION: UNITS, INTERCALIBRATION, INTEROPERABILITY, AND NORMALIZATION

826 The ability to make comparisons of results across global scales 827 of time and space in the oceans is a key appeal for embarking 828 on ocean metaproteomic research science. Indeed, the ability 829 of proteins to record the functional attributes of each 830 population of marine microbes could allow a "personalized 831 medicine" of the oceans, ²⁷ where long-term metaproteomic 832 records would track changes in environmental stresses 833 experienced by major taxa, and their resultant influence on

global biogeochemical cycles and implications for sustain- 834 ability. In ocean sciences, the few long-term time series 835 available allow studies of the impacts of global change on the 836 oceans. However, achieving the ambitious goal of integrating 837 metaproteome studies into global change science will require a 838 sufficient level of confidence regarding the accuracy and 839 precision of analyses to allow detection of changes between 840 samples sets. Metaproteomic data sets have reported 841 quantitative results in a variety of units thus far, including 842 total or normalized spectral counts, precursor intensities, or 843 calibrated absolute concentrations (fmol L^{-1} seawater). 844 Because the biological "matrix" of an environmental location 845 can change with time, there is a particular value in absolute 846 measurements that record peptide and protein abundances in 847 SI units per liter that can be unequivocally compared across 848 time. As a result, focusing on attributes that enable 849 interoperability between samples, even as technologies 850 (including chromatography, mass spectrometry, and infor- 851 matics) and reference databases improve, is an important 852 aspect of ocean metaproteomic data sharing. Efforts to 853 harmonize across analytical platforms to improve intercompar- 854 ability may be possible even in relative measurements 855 (nonabsolute) through the calibration of signal intensity 856 using a common reference material.⁹⁴ As described earlier, 857 efforts toward intercalibration of targeted metaproteomic 858 analyses, as well as intercomparison of global relative 859 abundance studies are critical to validating current measure- 860 ments and enabling future comparisons. Similarly, allowing 861 versioning of data sets will enable reanalysis of and 862 reinterpretation of historical data sets that can then be used 863 for temporal comparisons as both reference metagenomic 864 databases and PSM algorithms improve.

Another consideration in the use of metaproteomic data is 866 the choice of whether to normalize protein data to another 867 protein or parameter. This choice may reflect scientific culture 868 to some degree: in biological spheres normalization is routine 869 in order to provide organismal or ecological context, while in 870 chemical oceanography normalization is rarer due to an 871 appreciation for the importance of relaying absolute quantities 872 of molecules or elements per volume of seawater and the 873 fundamental interoperability of absolute units. Notably, 874 normalization approaches developed for single-organism 875 proteomics are not always applicable or appropriate for 876 metaproteomics. For example, assumptions of a constant 877 background proteome (in terms of uniformity of biological 878 organism(s) present) are not valid in many environments 879 depending on spatial or temporal scales being studied. 880 Moreover, the influence of differences in species abundance 881

882 across samples should be considered when considering 883 normalization of metaproteomic data sets; 95 for example, 884 when biological community composition changes across the 885 sampling regime normalization to a particular organism may 886 not be appropriate. While there have been advances in data 887 processing approaches that address aspects of this issue, 82 888 significant challenges remain.

PEPTIDE LEVEL REPORTING

890 Most publication guidelines for proteomics experiments 891 recommend at least two peptides be identified to confidently 892 report the identification of a specific protein. In the case of 893 metaproteomics, however, it is understood that SNPs, amino 894 acid variations, and substitutions associated with natural 895 biological diversity within species or strain-level populations 896 are common. As a result, it is possible to generate numerous 897 high-quality PSMs in metaproteomics that are "one-hit-898 wonders", likely due to a combination of challenges described 899 above, including limitations associated with the application of 900 mass spectrometry and PSM algorithms to highly complex 901 samples, availability of a suitable database due to limited 902 metagenomic coverage and quality, as well as the inherent 903 biological diversity of a given protein (and its host organism) 904 present in each nonclonal population found in the natural 905 environment. In most cases, the natural biological diversity 906 within species or strain-level populations is of great interest. It 907 is not uncommon for peptide sequences (typically tryptic 908 peptides) to be shared between closely related organisms. In 909 marine microbiology, it is becoming common for multiple 910 strains from a single species to have their genomes sequenced 911 and physiology studied. These strains are described as being 912 ecotypes that inhabit distinct environmental niches and can 913 have overlapping distributions allowing co-occurrence within 914 individual environmental samples. 96 As a result, the assignment 915 of multiple high-quality PSMs to a single protein sequence 916 derived from a single isolate genome sequence is not as 917 straightforward as with clonal populations of model laboratory 918 organisms. If multiple ecotypes are present with slight 919 variations in peptide sequence, the assignment of peptides to 920 protein sequences could need reconsideration. Indeed, this 921 subject intersects with the larger question regarding the 922 appropriate definition of microbial species itself. Nevertheless, 923 it has been demonstrated that this complexity can be taken 924 advantage of in order to design targeted metaproteomic 925 workflows that can interpret peptide biomarker abundances on 926 a large ocean biome scale and that multiple peptide biomarkers 927 provide consistent results.3

As a result, the two-peptide rule may not be appropriate for 929 metaproteomics. There is precedence for not using the two-930 peptide rule in splice variant analysis and detection of post-931 translationally modified amino acid sites. With the subsequent 932 arrival of high resolution mass spectrometry and stringent 933 FDR-based analyses, high-quality PSM that do not map to a 934 single protein sequence can have considerable value, and their 935 inability to have multiple peptides mapping to a protein are 936 likely related to numerous other challenges associated with 937 metaproteomic diversity and dynamic range as described 938 above, rather than necessarily being false positive identifica-939 tions. The ability to report multiple peptides to a specific 940 protein and its resultant percent peptide coverage itself 941 becomes associated with uncertainty if there are multiple 942 species with similar but not identical protein sequences. The 943 adoption of identification of protein families may become a

useful approach in metaproteomics where detection of 944 peptides with small variation in sequence diversity can 945 aggregate to a high confidence detection of a protein family 946 belonging to multiple ecotypes of a species, or a defined higher 947 taxonomic level particularly when interested in the bio- 948 geochemical impact of an enzyme.

Because of these challenges, targeted metaproteomic and 950 informatics efforts have focused on the tryptic peptide level, 951 using suites of tryptic peptide biomarkers as proxies for 952 proteins and processes of interest. Informatic tools such as 953 Unipept and Metatryp focus on tryptic peptides by conducting 954 analysis of shared tryptic peptides between different genomes 955 or metagenomes in order to maximize taxonomic interpreta-956 tion of peptide identification. While the consensus on what 957 should be considered the best practice for a high-quality 958 peptide identification is beyond the scope of this review, it is 959 clear that combining high-resolution mass spectrometry 960 capabilities, low false discovery rate, observance of peptides 961 in multiple spectra, visual inspection when possible, and other 962 factors can contribute to high-quality peptide identifications. 963

■ ENCOURAGING PROPER DATA USE

As metaproteomics is a relatively young data type, there is 965 potential for misunderstanding or misuse of results leading to 966 incorrect interpretations. These could have an inadvertent 967 detrimental effect of resulting in lost time chasing false leads or 968 loss of confidence in metaproteomic methods.⁹⁷ When used by 969 expert data generators, this risk is lessened due to a thorough 970 understanding of the limitations and methodology behind the 971 data. However, in the effort to share metaproteomic results 972 with a broader community of nonexpert users, there is 973 considerable risk that researchers will incorrectly attempt to 974 merge data units inappropriately and/or apply inappropriate 975 data transformations that could result in incorrect interpreta- 976 tion. For example, spectral counts are a popular quantitative 977 unit in proteomics that is powerful in assessing changes in each 978 individual protein's relative abundance across a range of 979 samples. However, efforts to compare abundances between 980 different proteins using relative abundance measurements such 981 as these should be minimized or replaced by calibrated 982 targeted measurements due to the variable influence of protein 983 size (and resultant number of tryptic peptides) and the 984 ionization efficiency of those peptides on each protein's 985 spectral count amplitude range. Nonexpert users may be 986 tempted to conduct meta-analyses of spectral count results that 987 could lead to faulty conclusions. As a result, efforts to educate 988 and encourage dialogue among data generators and nonexpert 989 users are important in fostering proper use of shared data sets. 990

Providing effective means for attribution of effort for those 991 involved in data generation is also important. Ideally, this could 992 include inviting the generator of a data set of interest to 993 collaborate and be a coauthor in studies. In addition, 994 acknowledging the use of a data set by citing original data 995 release manuscripts and DOI identifiers assigned to the data 996 set will be important in enabling data use to have metrics. The 997 attribution component is important in the sustainability of data 998 sharing projects, as this will incentivize the use of data sharing 999 portals and repositories by generators. If data generators feel 1000 they are not being properly attributed, they may be reluctant to 1001 share data and/or may seek more obscure avenues for meeting 1002 data sharing requirements. Learning about data use policy 1003 experience of prior metagenomics and large ocean programs 1004 such as GEOTRACES will be valuable in this regard.

Journal of Proteome Research

1091

1092

1100

1101

1106

1006 CONCLUSIONS

1007 Metaproteomics data sets have the potential to become a 1008 valuable data type to the ocean science community in that they 1009 represent a metabolic record of the status of the key microbial 1010 components within specific geographic environments through 1011 time. With significant regional and global ecosystem changes 1012 now occurring, 98 having access to detailed metabolic records 1013 through proteomic analyses of key environments could be 1014 particularly useful in providing an understanding of anthro-1015 pogenic impacts on natural ecosystems. Given that marine 1016 ecosystems are important to human society in a variety of 1017 ways, including maintaining Earth's habitability through 1018 microbial biogeochemical cycling, economic activities such as 1019 fisheries and aquaculture, and strategic and security importance 1020 to naval operations, the development and sharing of marine 1021 metaproteomic data sets will likely contribute to the long-term 1022 goal of developing a sustainable human society. This creates a 1023 distinct set of use cases for environmental proteomic data sets 1024 compared to those of laboratory cultivated organism or clinical 1025 proteomes, where planetary scale geospatial and temporal 1026 information are critically important metadata, and correspond-1027 ing environmental data are fundamental to contextualizing 1028 environmental metaproteomic results. Moreover, the amount 1029 of research funding going into biomedical proteomic research 1030 vastly outweighs comparable resources in the ocean environ-1031 ment, making scarce ocean data sets of considerable value. 1032 These underlying differences in data usage and investment 1033 between environmental and biomedical proteomic data sets 1034 demonstrate a need for distinct data sharing strategies, and we 1035 have proposed some best practices with regard to metadata 1036 needs for ocean metaproteome data sharing, as well as 1037 summarized challenges associated with conducting metapro-1038 teomic research in hopes of inspiring innovation and 1039 collaboration.

1040 EXPERIMENTAL METHODS

1041 While the data presented in this review manuscript were largely 1042 previously published,³ some novel interpretations of the data 1043 have been included to demonstrate the complexity of 1044 metaproteome samples. The methods for these comparisons 1045 are briefly described below. A human cell line (HeLa) and 1046 ocean metaproteome samples (METZYME KM1128, Station 5 1047 0°N 158°W 40, 60, 120, 150, and 200 m depth, 0.2 μ m filter 1048 pore size, prefiltered with 3.0 μ m pore size) were analyzed 1049 under identical chromatographic and mass spectrometry 1050 conditions to provide examples of sample complexity run 1051 within 5 days of each other. Protein extraction for 1052 metaproteomics was conducted using SDS detergent and 1053 tube gel purification as previously described.³ Protein extracts 1054 were analyzed by liquid chromatography-mass spectrometry 1055 (LC-MS) (Michrom Advance HPLC coupled to a Thermo 1056 Scientific Fusion Orbitrap mass spectrometer with a Thermo 1057 Flex source). A total of 0.5 μg (HeLa) or 1 μg (ocean) of each 1058 sample (measured before trypsin digestion) was concentrated 1059 onto a trap column (0.2 × 10 mm ID, 5 μ m particle size, 120 Å 1060 pore size, C18 Reprosil-Gold, Dr. Maisch GmbH) and rinsed 1061 with 100 μ L of 0.1% formic acid, 2% acetonitrile (ACN), 1062 97.9% water before gradient elution through a reverse phase 1063 C18 nanospray column (0.1 \times 400 mm ID, 3 μ m particle size, 1064 120 Å pore size, C18 Reprosil-Gold, Dr. Maisch GmbH) at a 1065 flow rate of 300 nL/min. The chromatography consisted of a 1066 nonlinear 200 min gradient from 5% to 95% buffer B, where A

was 0.1% formic acid in water and B was 0.1% formic acid in 1067 ACN (all solvents were Fisher Optima grade). The mass 1068 spectrometer was set to perform MS scans on the Orbitrap 1069 (240 000 resolution at 200 m/z) with a scan range of 380 m/z 1070 to 1580 m/z. MS/MS was performed on the ion trap using 1071 data-dependent settings (top speed, dynamic exclusion 15 s, 1072 excluding unassigned and singly charged ions, precursor mass 1073 tolerance of ± 3 ppm, with a maximum injection time of 150 1074 ms).

Quantitive Peak Comparisons

Comparisons of the relative complexity of ocean metaproteo- 1077 mic samples (METZYME expedition, Station 3, depths 40, 60, 1078 120, 150, and 250 m) with a human cell line sample run in 1079 triplicate was conducted. A total of 0.5 μg of Hela was injected 1080 per replicate, while 1 μg of ocean metaproteomic sample was 1081 mjected per sample, as described above. Precursor peak data 1082 were extracted from raw files using ProteoWizard's MSCon- 1083 vertGUI to text file using the vendor (Thermo) peak picking 1084 algorithm, and applying two filters: the peak picking algorithm 1085 (set for MS Levels 1 only) followed by MS level filter MS level 1086 1 only. The number of precursor peaks per MS1, total ion 1087 count (TIC), and chromatographic time information were 1088 then extracted from the output files using a custom Python 1089 script, summed, and visualized (Figure 5).

ASSOCIATED CONTENT

Supporting Information

The Supporting Information is available free of charge on the 1093 ACS Publications website at DOI: 10.1021/acs.jproteo- 1094 me.8b00761.

Table S1. Examples of environmental and biomedical 1096 metaproteomics studies. Table S2. Additional metadata 1097 associated with sample processing and analysis for, but 1098 not exclusive to, metaproteomics data (PDF) 1099

AUTHOR INFORMATION

Corresponding Author

*Address: Marine Chemistry and Geochemistry Department, 1102 Woods Hole Oceanographic Institution, 360 Woods Hole 1103 Road, Woods Hole, MA 02543. E-mail: msaito@whoi.edu. 1104 Phone: 1-508-289-2393.

ORCID ®

Mak A. Saito: 0000-0001-6040-9295	1107
Megan E. Duffy: 0000-0002-3212-4927	1108
David A. Gaylord: 0000-0001-7987-6870	1109
Noelle A. Held: 0000-0003-1073-0851	1110
William Judson Hervey, IV: 0000-0003-3285-6754	1111
Robert L. Hettich: 0000-0001-7708-786X	1112
Pratik Jagtap: 0000-0003-0984-0973	1113
Michael G. Janech: 0000-0002-3202-4811	1114
Danie B. Kinkade: 0000-0002-1134-7347	1115
Dasha Leary: 0000-0003-2325-7143	1116
Matthew McIlvin: 0000-0002-5301-8365	1117
Eli Moore: 0000-0002-9750-7769	1118
Benjamin A. Neely: 0000-0001-6120-7695	1119
Jaclyn K. Saunders: 0000-0003-1023-6239	1120
Adam Shepherd: 0000-0003-4486-9448	1121
Nicholas Symmonds: 0000-0002-9436-0351	1122
David Walsh: 0000-0002-9951-5447	1123

1124 Notes

1125 The authors declare no competing financial interest.

1126 The data in this manuscript was previously published^{3,14} and is

- 1127 available at the Biological and Chemical Oceanography Data
- 1128 Management Office repository (https://www.bco-dmo.org/
- 1129 dataset/646115). The raw mass spectra files are available at
- 1130 PRIDE under project name Pacific Ocean metaproteomics
- 1131 METZYME KM1128, Project accession: PXD009712; Project
- 1132 DOI: 10.6019/PXD009712.

1133 **ACKNOWLEDGMENTS**

1134 The workshop that led to this manuscript was funded by an 1135 NSF EarthCube Grant No. 1639714 and an anonymous donor 1136 to M.A.S. and D.K. We thank Mary Zawoysky for assistance 1137 with workshop planning and manuscript editing and Natalie 1138 Renier for graphics assistance.

REFERENCES

- (1) Sowell, S. M.; Wilhelm, L. J.; Norbeck, A. D.; Lipton, M. S.; 1141 Nicora, C. D.; Barofsky, D. F.; Carlson, C. A.; Smith, R. D.; 1142 Giovanonni, S. J. Transport functions dominate the SAR11 1143 metaproteome at low-nutrient extremes in the Sargasso Sea. ISME
- 1144 *J.* **2009**, 3, 93–105.

1184 1431-1450.

1188 ogr.: Methods 2012, 10 (5), 353-366.

- (2) Morris, R. M.; Nunn, B. L.; Frazar, C.; Goodlett, D. R.; Ting, Y. 1146 S.; Rocap, G. Comparative metaproteomics reveals ocean-scale shifts 1147 in microbial nutrient utilization and energy transduction. ISME J. 1148 **2010**, *4*, 673–685.
- (3) Saito, M. A.; McIlvin, M. R.; Moran, D. M.; Goepfert, T. J.; 1150 DiTullio, G. R.; Post, A. F.; Lamborg, C. H. Multiple nutrient stresses 1151 at intersecting Pacific Ocean biomes detected by protein biomarkers. 1152 Science 2014, 345 (6201), 1173-1177.
- (4) Hawley, A. K.; Brewer, H. M.; Norbeck, A. D.; Paša-Tolić, L.; 1154 Hallam, S. J. Metaproteomics reveals differential modes of metabolic 1155 coupling among ubiquitous oxygen minimum zone microbes. Proc. 1156 Natl. Acad. Sci. U. S. A. 2014, 111 (31), 11395-11400.
- 1157 (5) Leary, D. H.; Li, R. W.; Hamdan, L. J.; Hervey, W. J., IV; 1158 Lebedev, N.; Wang, Z.; Deschamps, J. R.; Kusterbeck, A. W.; Vora, G. 1159 J. Integrated metagenomic and metaproteomic analyses of marine 1160 biofilm communities. *Biofouling* **2014**, 30 (10), 1211–1223.
- (6) Moore, E. K.; Harvey, H. R.; Faux, J. F.; Goodlett, D. R.; Nunn, 1162 B. L. Protein recycling in Bering Sea algal incubations. Mar. Ecol.: 1163 Prog. Ser. 2014, 515, 45-59.
- (7) Moore, E. K.; Nunn, B. L.; Goodlett, D. R.; Harvey, H. R. 1165 Identifying and tracking proteins through the marine water column: 1166 Insights into the inputs and preservation mechanisms of protein in 1167 sediments. Geochim. Cosmochim. Acta 2012, 83, 324-359.
- (8) Bridoux, M. C.; Neibauer, J.; Ingalls, A. E.; Nunn, B. L.; Keil, R. 1169 G. Suspended marine particulate proteins in coastal and oligotrophic 1170 waters. Journal of Marine Systems 2015, 143, 39-48.
- (9) Georges, A. A.; El-Swais, H.; Craig, S. E.; Li, W. K.; Walsh, D. A. 1172 Metaproteomic analysis of a winter to spring succession in coastal 1173 northwest Atlantic Ocean microbial plankton. ISME J. 2014, 8 (6), 1174 1301.
- (10) Bender, S. J.; Moran, D. M.; McIlvin, M. R.; Zheng, H.; 1175 1176 McCrow, J. P.; Badger, J.; DiTullio, G. R.; Allen, A. E.; Saito, M. A. 1177 Colony formation in Phaeocystis antarctica connecting molecular 1178 mechanisms with iron biogeochemistry. Biogeosciences 2018, 15 (16), 1179 4923-4942.
- (11) Bertrand, E. M.; Moran, D. M.; McIlvin, M. R.; Hoffman, J. M.; 1181 Allen, A. E.; Saito, M. A. Methionine synthase interreplacement in 1182 diatom cultures and communities: Implications for the persistence of 1183 B₁₂ use by eukaryotic phytoplankton. Limnol. Oceanogr. 2013, 58 (4),
- 1185 (12) Moore, E. K.; Nunn, B. L.; Faux, J. F.; Goodlett, D. R.; Harvey, 1186 H. R. Evaluation of electrophoretic protein extraction and database-1187 driven protein identification from marine sediments. Limnol. Ocean-

- (13) Kan, J.; Hanson, T.; Ginter, J.; Wang, K.; Chen, F. 1189 Metaproteomic analysis of Chesapeake Bay microbial communities. 1190 Saline Syst. 2005, 1 (1), 7.
- (14) Saito, M. A.; Dorsk, A.; Post, A. F.; McIlvin, M.; Rappé, M. S.; 1192 DiTullio, G.; Moran, D. Needles in the Blue Sea: Sub Species 1193 Specificity in Targeted Protein Biomarker Analyses Within the Vast 1194 Oceanic Microbial Metaproteome. Proteomics 2015, 15 (20), 3521-1195
- (15) Williams, T. J.; Long, E.; Evans, F.; DeMaere, M. Z.; Lauro, F. 1197 M.; Raftery, M. J.; Ducklow, H.; Grzymski, J. J.; Murray, A. E.; 1198 Cavicchioli, R. A metaproteomic assessment of winter and summer 1199 bacterioplankton from Antarctic Peninsula coastal surface waters. 1200 ISME J. 2012, 6 (10), 1883-1900.
- (16) Bertrand, E. M.; Allen, A. E.; Dupont, C. L.; Norden-Krichmar, 1202 T.; Bai, J.; Saito, M. A.; Valas, R. E. Influence of Cobalamin Starvation 1203 on Diatom Molecular Physiology and the Identification of a 1204 Cobalamin Acquisition Protein. Proc. Natl. Acad. Sci. U. S. A. 2012, 1205 109 (26), E1762-E1771.
- (17) Dyhrman, S. T.; Jenkins, B. D.; Rynearson, T. A.; Saito, M. A.; 1207 Mercier, M. L.; Alexander, H.; Whitney, L. P.; Drzewianowski, A.; 1208 Bulygin, V. V.; Bertrand, E. M.; Wu, Z.; Benitez-Nelson, C.; Heithoff, 1209 A. The transcriptome and proteome of the diatom Thalassiosira 1210 pseudonana reveal a diverse phosphorus stress response. PLoS One 1211 **2012**, 7 (3), No. e33768.
- (18) Cox, A. D.; Saito, M. A. Proteomic responses of oceanic 1213 Synechococcus WH8102 to phosphate and zinc scarcity and cadmium 1214 additions. Front. Microbiol. 2013, 4, 387.
- (19) Mackey, K. R.; Post, A. F.; McIlvin, M. R.; Cutter, G. A.; John, 1216 S. G.; Saito, M. A. Divergent responses of Atlantic coastal and oceanic 1217 Synechococcus to iron limitation. Proc. Natl. Acad. Sci. U. S. A. 2015, 1218 112 (32), 9944-9949.
- (20) Wurch, L. L.; Bertrand, E. M.; Saito, M. A.; Van Mooy, B. A. S.; 1220 Dyhrman, S. T. Proteome Changes Driven by Phosphorus Deficiency 1221 and Recovery in the Brown Tide-Forming Alga Aureococcus 1222 anophagefferens. PLoS One 2011, 6 (12), No. e28949.
- (21) Swanner, E. D.; Wu, W.; Hao, L.; Wüstner, M. L.; Obst, M.; 1224 Moran, D. M.; McIlvin, M. R.; Saito, M. A.; Kappler, A. Physiology, 1225 Fe (II) oxidation, and Fe mineral formation by a marine planktonic 1226 cyanobacterium grown under ferruginous conditions. Frontiers in 1227 Earth Science 2015, 3, 60.
- (22) Nunn, B.; Aker, J.; Shaffer, S.; Tsai, Y.; Strzepek, R.; Boyd, P.; 1229 Freeman, T.; Brittnacher, M.; Malmstrom, L.; Goodlett, D. Decipher- 1230 ing diatom biochemical pathways via whole-cell proteomics. Aquat. 1231 Microb. Ecol. 2009, 55 (3), 241-253.
- (23) Nunn, B. L.; Slattery, K. V.; Cameron, K. A.; Timmins- 1233 Schiffman, E.; Junge, K. Proteomics of Colwellia psychrerythraea at 1234 subzero temperatures-a life with limited movement, flexible 1235 membranes and vital DNA repair. Environ. Microbiol. 2015, 17 (7), 1236
- (24) Heyer, R.; Schallert, K.; Zoun, R.; Becher, B.; Saake, G.; 1238 Benndorf, D. Challenges and perspectives of metaproteomic data 1239 analysis. J. Biotechnol. 2017, 261, 24-36.
- (25) Saito, M. A.; Bulygin, V. V.; Moran, D. M.; Taylor, C.; Scholin, 1241 C. Examination of Microbial Proteome Preservation Techniques 1242 Applicable to Autonomous Environmental Sample Collection. Front. 1243 Microbiol. 2011, 2, 215.
- (26) Breier, J.; Gomez-Ibanez, D.; Reddington, E.; Huber, J.; 1245 Emerson, D. A precision multi-sampler for deep-sea hydrothermal 1246 microbial mat studies. Deep Sea Res., Part I 2012, 70, 83-90.
- (27) Saito, M. A.; Breier, C.; Jakuba, M.; McIlvin, M.; Moran, D. 1248 Envisioning a Chemical Metaproteomics Capability for Biochemical 1249 Research and Diagnosis of global Ocean Microbiomes. In The 1250 Chemistry of Microbiomes: Proceedings of a Seminar Series; National 1251 Academies Press: National Academies of Sciences, Engineering, 1252 Medicine, 2017; pp 29–36.
- (28) Boyle, E. A.; Bergquist, B. A.; Kayser, R. A.; Mahowald, N. Iron, 1254 manganese, and lead at Hawaii Ocean Time-series station ALOHA: 1255 Temporal variability and an intermediate water hydrothermal plume. 1256 Geochim. Cosmochim. Acta 2005, 69 (4), 933-952.

- 1258 (29) Bergauer, K.; Fernandez-Guerra, A.; Garcia, J. A. L.; Sprenger, 1259 R. R.; Stepanauskas, R.; Pachiadaki, M. G.; Jensen, O. N.; Herndl, G. 1260 J. Organic matter processing by microbial communities throughout 1261 the Atlantic water column as revealed by metaproteomics. *Proc. Natl.* 1262 *Acad. Sci. U. S. A.* **2018**, *115* (3), E400–408.
- 1263 (30) Venter, J. C.; Remington, K.; Heidelberg, J. F.; Halpern, A. L.; 1264 Rusch, D.; Eisen, J. A.; Wu, D.; Paulsen, I.; Nelson, K. E.; Nelson, W.; 1265 Fouts, D. E.; Levy, S.; Knap, A. H.; Lomas, M. W.; Nealson, K.; 1266 White, O.; Peterson, J.; Hoffman, J.; Parsons, R.; Baden-Tillson, H.; 1267 Pfannkoch, C.; Rogers, Y. H.; Smith, H. O. Environmental genome 1268 shotgun sequencing of the Sargasso Sea. *Science* 2004, 304 (5667), 1269 66–74.
- 1270 (31) DeLong, E. F.; Preston, C. M.; Mincer, T.; Rich, V.; Hallam, S. 1271 J.; Frigaard, N.-U.; Martinez, A.; Sullivan, M. B.; Edwards, R.; Brito, B. 1272 R.; Chisholm, S. W.; Karl, D. M. Community Genomics Among 1273 Stratified Microbial Assemblages in the Ocean's Interior. *Science* 2006, 1274 311 (5760), 496–503.
- 1275 (32) Nunn, B. L.; Keil, R. G. A comparison of non-hydrolytic 1276 methods for extracting amino acids and proteins from coastal marine 1277 sediments. *Mar. Chem.* **2006**, *98* (1), 31–42.
- 1278 (33) Keil, R. G.; Tsamakis, E.; Fuh, C. B.; Giddings, J. C.; Hedges, J. 1279 I. Mineralogical and textural controls on the organic composition of 1280 coastal marine sediments: hydrodynamic separation using SPLITT-1281 fractionation. *Geochim. Cosmochim. Acta* 1994, 58 (2), 879–893.
- 1282 (34) Ram, R. J.; VerBerkmoes, N. C.; Thelen, M. P.; Tyson, G. W.; 1283 Baker, B. J.; Blake, R. C., II; Shah, M.; Hettich, R. L.; Banfield, J. F. 1284 Community Proteomics of a Natural Microbial Biofilm. *Science* 2005, 1285 308 (5730), 1915–1920.
- 1286 (35) VerBerkmoes, N. C.; Hervey, W. J.; Shah, M.; Land, M.; 1287 Hauser, L.; Larimer, F. W.; Van Berkel, G. J.; Goeringer, D. E. 1288 Evaluation of "shotgun" proteomics for identification of biological 1289 threat agents in complex environmental matrixes: experimental 1290 simulations. *Anal. Chem.* 2005, 77 (3), 923–932.
- 1291 (36) Mattes, T. E.; Nunn, B. L.; Marshall, K. T.; Proskurowski, G.; 1292 Kelley, D. S.; Kawka, O. E.; Goodlett, D. R.; Hansell, D. A.; Morris, R. 1293 M. Sulfur oxidizers dominate carbon fixation at a biogeochemical hot 1294 spot in the dark ocean. *ISME J.* 2013, 7 (12), 2349.
- 1295 (37) Sinitcyn, P.; Rudolph, J. D.; Cox, J. Computational Methods for 1296 Understanding Mass Spectrometry—Based Shotgun Proteomics Data. 1297 Annual Review of Biomedical Data Science 2018, 1 (1), 207–234.
- 1298 (38) Gillet, L. C.; Navarro, P.; Tate, S.; Röst, H.; Selevsek, N.; 1299 Reiter, L.; Bonner, R.; Aebersold, R. Targeted data extraction of the 1300 MS/MS spectra generated by data-independent acquisition: a new 1301 concept for consistent and accurate proteome analysis. *Mol. Cell.* 1302 *Proteomics* **2012**, *11* (6), O111–O111.016717.
- 1303 (39) Gallien, S.; Bourmaud, A.; Kim, S. Y.; Domon, B. Technical 1304 considerations for large-scale parallel reaction monitoring analysis. *J.* 1305 *Proteomics* **2014**, *100*, 147–159.
- 1306 (40) Gallien, S.; Duriez, E.; Crone, C.; Kellmann, M.; Moehring, T.; 1307 Domon, B. Targeted proteomic quantification on quadrupole-orbitrap 1308 mass spectrometer. *Mol. Cell. Proteomics* **2012**, *11* (12), 1709–1723. 1309 (41) Gallien, S.; Duriez, E.; Domon, B. Selected reaction monitoring 1310 applied to proteomics. *J. Mass Spectrom.* **2011**, *46* (3), 298–312.
- 1311 (42) Pluskal, T.; Castillo, S.; Villar-Briones, A.; Orešič, M. MZmine 1312 2: Modular framework for processing, visualizing, and analyzing mass 1313 spectrometry-based molecular profile data. *BMC Bioinf.* **2010**, *11* (1), 1314 395.
- 1315 (43) Nunn, B. L.; Faux, J. F.; Hippmann, A. A.; Maldonado, M. T.; 1316 Harvey, H. R.; Goodlett, D. R.; Boyd, P. W.; Strzepek, R. F. Diatom 1317 proteomics reveals unique acclimation strategies to mitigate Fe 1318 limitation. *PLoS One* **2013**, 8 (10), No. e75653.
- 1319 (44) Dorfer, V.; Maltsev, S.; Winkler, S.; Mechtler, K. CharmeRT: 1320 Boosting peptide identifications by chimeric spectra identification and 1321 retention time prediction. *J. Proteome Res.* **2018**, *17* (8), 2581–2589. 1322 (45) Cutter, G. A. Intercalibration in chemical oceanography—1323 getting the right number. *Limnol. Oceanogr.: Methods* **2013**, *11* (7), 1324 418–424.

- (46) The, M.; Tasnim, A.; Käll, L. How to talk about protein-level 1325 false discovery rates in shotgun proteomics. *Proteomics* **2016**, *16* (18), 1326 2461–2469.
- (47) Noble, W. S. Mass spectrometrists should search only for 1328 peptides they care about. *Nat. Methods* **2015**, *12* (7), 605–608.
- (48) Huang, T.; Wang, J. J.; Yu, W. C.; He, Z. Protein inference: a 1330 review. *Briefings Bioinf.* **2012**, *13* (5), 586–614.
- (49) Nesvizhskii, A. I.; Aebersold, R. Interpretation of shotgun 1332 proteomic data The protein inference problem. *Mol. Cell. Proteomics* 1333 **2005**, *4* (10), 1419–1440.
- (50) Brosch, M.; Saunders, G. I.; Frankish, A.; Collins, M. O.; et al. 1335 Shotgun proteomics aids discovery of novel protein coding genes, 1336 alternative splicing, and "resurrected" pseudogenes in the mouse 1337 genome. *Genome Res.* **2011**, 21 (5), 756–767.
- (51) Volders, P. J.; Verheggen, K.; Menschaert, G.; Vandepoele, K.; 1339 et al. An update on LNCipedia: a database for annotated human 1340 lncRNA sequences. *Nucleic Acids Res.* **2015**, 43 (8), 4363–4364. 1341
- (52) Slavoff, S. A.; Mitchell, A. J.; Schwaid, A. G.; Cabili, M. N.; Ma, 1342 J.; Levin, J. Z.; Karger, A. D.; Budnik, B. A.; Rinn, J. L.; Saghatelian, A. 1343 Peptidomic discovery of short open reading frame-encoded peptides 1344 in human cells. *Nat. Chem. Biol.* **2013**, *9*, 59–64.
- (53) Muth, T.; Kolmeder, C. A.; Salojarvi, J.; Keskitalo, S.; Varjosalo, 1346 M.; Verdam, F. J.; Rensen, S. S.; Reichl, U.; de Vos, W. M.; Rapp, E.; 1347 Martens, L. Navigating through metaproteomics data: A logbook of 1348 database searching. *Proteomics* **2015**, *15* (20), 3439–3453.
- (54) Timmins-Schiffman, E.; May, D. H.; Mikan, M.; Riffle, M.; 1350 Frazar, C.; Harvey, H.; Noble, W. S.; Nunn, B. L. Critical decisions in 1351 metaproteomics: achieving high confidence protein annotations in a 1352 sea of unknowns. *ISME J.* **2017**, *11* (2), 309–314.
- (55) Teeling, H.; Fuchs, B. M.; Becher, D.; Klockow, C.; 1354 Gardebrecht, A.; Bennke, C. M.; Kassabgy, M.; Huang, S.; Mann, 1355 A. J.; Waldmann, J.; Weber, M.; Klindworth; Otto, A.; Lange, J.; 1356 Bernhardt, J.; Reinsch, C.; Hecker, M.; Peplies, J.; Bockelmann, F. D.; 1357 Callies, U.; Gerdts, G.; Wichels, A.; Wiltshire, K. H.; Glockner, F. O.; 1358 Schweder, T.; Amann, R. Substrate-controlled succession of marine 1359 bacterioplankton populations induced by a phytoplankton bloom. 1360 Science 2012, 336 (6081), 608–611.
- (56) Colatriano, D.; Ramachandran, A.; Yergeau, E.; Maranger, R.; 1362 Gélinas, Y.; Walsh, D. A. Metaproteomics of aquatic microbial 1363 communities in a deep and stratified estuary. *Proteomics* **2015**, *15* 1364 (20), 3566–3579.
- (57) Burton, J. N.; Liachko, I.; Dunham, M. J.; Shendure, J. Species- 1366 level deconvolution of metagenome assemblies with Hi-C-based 1367 contact probability maps. *G3: Genes, Genomes, Genet.* **2014**, *4*, 1339– 1368 1346.
- (58) Kashtan, N.; Roggensack, S. E.; Rodrigue, S.; Thompson, J. W.; 1370 Biller, S. J.; Coe, A.; Ding, H.; Marttinen, P.; Malmstrom, R. R.; 1371 Stocker, R.; Follows, M. J.; Stepanauskas, R.; Chisholm, S. W. Single-1372 Cell Genomics Reveals Hundreds of Coexisting Subpopulations in 1373 Wild Prochlorococcus. *Science* 2014, 344 (6182), 416–420.
- (59) Rusch, D. B.; Martiny, A. C.; Dupont, C. L.; Halpern, A. L.; 1375 Venter, J. C. Characterization of *Prochlorococcus* clades from iron- 1376 depleted oceanic regions. *Proc. Natl. Acad. Sci. U. S. A.* **2010**, 107 1377 (37), 16184–16189.
- (60) Yooseph, S.; Sutton, G.; Rusch, D. B.; Halpern, A. L.; 1379 Williamson, S. J.; Remington, K.; Eisen, J. A.; Heidelberg, K. B.; 1380 Manning, G.; Li, W.; Jaroszewski, L.; Cieplak, P.; Miller, C. S.; Li, H.; 1381 Mashiyama, S. T.; Joachimiak, M. P.; van Belle, C.; Chandonia, J.-M.; 1382 Soergel, D. A.; Zhai, Y.; Natarajan, K.; Lee, S.; Raphael, B. J.; Bafna, 1383 V.; Friedman, R.; Brenner, S. E.; Godzik, A.; Eisenberg, D.; Dixon, J. 1384 E.; Taylor, S. S.; Strausberg, R. L.; Frazier, M.; Venter, J. C. The 1385 Sorcerer II Global Ocean Sampling Expedition: Expanding the 1386 Universe of Protein Families. *PLoS Biol.* 2007, 5 (3), No. e16.
- (61) Williamson, A. J.; Smith, D. L.; Blinco, D.; Unwin, R. D.; 1388 Pearson, S.; Wilson, C.; Miller, C.; Lancashire, L.; Lacaud, G.; 1389 Kouskoff, V.; Whetton, A. D. Quantitative proteomics analysis 1390 demonstrates post-transcriptional regulation of embryonic stem cell 1391 differentiation to hematopoiesis. *Mol. Cell. Proteomics* **2008**, 7 (3), 1392 459–472.

- 1394 (62) Keeling, P. J.; Burki, F.; Wilcox, H. M.; Allam, B.; Allen, E. E.; 1395 Amaral-Zettler, L. A.; Armbrust, E. V.; Archibald, J. M.; Bharti, A. K.; 1396 Bell, C. J.; Beszteri, B. The Marine Microbial Eukaryote Tran-1397 scriptome Sequencing Project (MMETSP): illuminating the func-1398 tional diversity of eukaryotic life in the oceans through transcriptome 1399 sequencing. *PLoS Biology* **2014**, *12* (6), No. e1001889.
- 1400 (63) May, D. H.; Timmins-Schiffman, E.; Mikan, M. P.; Harvey, H. 1401 R.; Borenstein, E.; Nunn, B. L.; Noble, W. S. An alignment-free 1402 "metapeptide" strategy for metaproteomic characterization of micro-1403 biome samples using shotgun metagenomic sequencing. *J. Proteome* 1404 *Res.* **2016**, 15 (8), 2697–2705.
- 1405 (64) Biller, S. J.; Berube, P. M.; Berta-Thompson, J. W.; Kelly, L.; 1406 Roggensack, S. E.; Awad, L.; Roache-Johnson, K. H.; Ding, H.; 1407 Giovannoni, S. J.; Rocap, G.; Moore, L. R.; Chisholm, S. W. Genomes 1408 of diverse isolates of the marine cyanobacterium *Prochlorococcus. Sci.* 1409 *Data* **2014**, *1*, 140034.
- 1410 (65) Grote, J.; Thrash, J. C.; Huggett, M. J.; Landry, Z. C.; Carini, 1411 P.; Giovannoni, S. J.; Rappé, M. S. Streamlining and core genome 1412 conservation among highly divergent members of the SAR11 clade. 1413 *mBio* **2012**, *3* (5), e00252–12.
- 1414 (66) Santoro, A. E.; Dupont, C. L.; Richter, R. A.; Craig, M. T.; 1415 Carini, P.; McIlvin, M. R.; Yang, Y.; Orsi, W. D.; Moran, D. M.; Saito, 1416 M. A. Genomic and proteomic characterization of "Candidatus 1417 Nitrosopelagicus brevis": an ammonia-oxidizing archaeon from the 1418 open ocean. Proc. Natl. Acad. Sci. U. S. A. 2015, 112 (4), 1173–1178. 1419 (67) Eng, J.; McCormack, A.; Yates, J. An approach to correlate 1420 tandem mass spectral data of peptides with amino acid sequences in a 1421 protein database. J. Am. Soc. Mass Spectrom. 1994, 5 (11), 976–989. 1422 (68) Eng, J. K.; Jahan, T. A.; Hoopmann, M. R. Comet: an open 1423 source MS/MS sequence database search tool. Proteomics 2013, 13 1424 (1), 22–24.
- 1425 (69) Craig, R.; Beavis, R. C. A method for reducing the time 1426 required to match protein sequences with tandem mass spectra. *Rapid* 1427 *Commun. Mass Spectrom.* **2003**, *17* (20), 2310–2316.
- 1428 (70) Jagtap, P.; McGowan, T.; Bandhakavi, S.; Tu, Z. J.; Seymour, 1429 S.; Griffin, T. J.; Rudney, J. D. Deep metaproteomic analysis of human 1430 salivary supernatant. *Proteomics* **2012**, *12* (7), 992–1001.
- 1431 (71) Jagtap, P.; Goslinga, J.; Kooren, J. A.; McGowan, T.; 1432 Wroblewski, M. S.; Seymour, S. L.; Griffin, T. J. A two step database 1433 search method improves sensitivity in peptide sequence matches for 1434 metaproteomics and proteogenomics studies. *Proteomics* **2013**, *13* (8), 1435 1352–1357.
- 1436 (72) Jagtap, P. D.; Johnson, J. E.; Onsongo, G.; Sadler, F. W.; 1437 Murray, K.; Wang, Y.; Shenykman, G. M.; Bandhakavi, S.; Smith, L. 1438 M.; Griffin, T. J. Flexible and accessible workflows for improved 1439 proteogenomic analysis using the Galaxy framework. *J. Proteome Res.* 1440 **2014**, *13* (12), 5898–5908.
- 1441 (73) Jagtap, P. D.; Blakely, A.; Murray, K.; Stewart, S.; Kooren, J.; 1442 Johnson, J. E.; Rhodus, N. L.; Rudney, J.; Griffin, T. J. Metaproteomic 1443 analysis using the Galaxy framework. *Proteomics* **2015**, *15* (20), 1444 3553–3565.
- 1445 (74) Rudney, J. D.; Jagtap, P. D.; Reilly, C. S.; Chen, R.; Markowski, 1446 T. W.; Higgins, L.; Johnson, J. E.; Griffin, T. J. Protein relative 1447 abundance patterns associated with sucrose-induced dysbiosis are 1448 conserved across taxonomically diverse oral microcosm biofilm 1449 models of dental caries. *Microbiome* 2015, 3 (1), 69.
- 1450 (75) Kertesz-Farkas, A.; Keich, U.; Noble, W. S. Tandem mass 1451 spectrum identification via cascaded search. *J. Proteome Res.* **2015**, *14* 1452 (8), 3027–3038.
- 1453 (76) Vaudel, M.; Burkhart, J. M.; Zahedi, R. P.; Oveland, E.; Berven, 1454 F. S.; Sickmann, A.; Martens, L.; Barsnes, H. PeptideShaker enables 1455 reanalysis of MS-derived proteomics data sets. *Nat. Biotechnol.* **2015**, 1456 33, 22.
- 1457 (77) Muth, T.; Hartkopf, F.; Vaudel, M.; Renard, B. Y. A Potential 1458 Golden Age to Come—Current Tools, Recent Use Cases, and Future 1459 Avenues for De Novo Sequencing in Proteomics. *Proteomics* **2018**, *18* 1460 (18), 1700150.
- 1461 (78) Mesuere, B.; van der Jeugt, F.; Devreese, B.; Vandamme, P.; 1462 Dawyndt, P. The unique peptidome: Taxon-specific tryptic peptides

- as biomarkers for targeted metaproteomics. *Proteomics* **2016**, *16* (17), 1463 2313–2318.
- (79) Mesuere, B.; Devreese, B.; Debyser, G.; Aerts, M.; Vandamme, 1465 P.; Dawyndt, P. Unipept: Tryptic Peptide-Based Biodiversity Analysis 1466 of Metaproteome Samples. *J. Proteome Res.* **2012**, *11* (12), 5773–1467 5780
- (80) Mesuere, B.; Willems, T.; van der Jeugt, F.; Devreese, B.; 1469 Vandamme, P.; Dawyndt, P. Unipept web services for metaproteo- 1470 mics analysis. *Bioinformatics* **2016**, *32* (11), 1746–1748.
- (81) Muth, T.; Kohrs, F.; Heyer, R.; Benndorf, D.; Rapp, E.; Reichl, 1472 U.; Martens, L.; Renard, B. Y. MPA Portable: A Stand-Alone Software 1473 Package for Analyzing Metaproteome Samples on the Go. *Anal. Chem.* 1474 **2018**, 90 (1), 685–689.
- (82) Riffle, M.; May, D.; Timmins-Schiffman, E.; Mikan, M. P.; 1476 Jaschob, D.; Noble, W. S.; Nunn, B. L. MetaGOmics: A Web-Based 1477 Tool for Peptide Centric Functional and Taxonomic Analysis of 1478 Metaproteomic Data. *Proteomes* **2018**, *6* (1), 2. 1479
- (83) Gurdeep Singh, R.; Tanca, A.; Palomba, A.; van der Jeugt, F.; 1480 Verschaffelt, P.; Uzzau, S.; Martens, L.; Dawyndt, P.; Mesuere, B. 1481 Unipept 4.0: Functional Analysis of Metaproteome Data. *J. Proteome* 1482 *Res.* 2019, 18, 606–615.
- (84) Huson, D. H.; Auch, A. F.; Qi, J.; Schuster, S. C. MEGAN 1484 analysis of metagenomic data. *Genome Res.* **2007**, 17 (3), 377–386. 1485
- (85) Huerta-Cepas, J.; Forslund, K.; Coelho, L. P.; Szklarczyk, D.; 1486 Jensen, L. J.; von Mering, C.; Bork, P. Fast Genome-Wide Functional 1487 Annotation through Orthology Assignment by eggNOG-Mapper. 1488 Mol. Biol. Evol. 2017, 34 (8), 2115–2122. 1489
- (86) Schnoes, A. M.; Brown, S. D.; et al. Annotation error in public 1490 databases: misannotation of molecular function in enzyme super- 1491 families. *PLoS Comput. Biol.* **2009**, *5* (12), No. e1000605.
- (87) Brenner, S. E. Errors in genome annotation. *Trends Genet.* 1493 **1999**, *15* (4), 132–133.
- (88) McQuaid, J. B.; Kustka, A. B.; Oborník, M.; Horák, A.; 1495 McCrow, J. P.; Karas, B. J.; Zheng, H.; Kindeberg, T.; Andersson, A. 1496 J.; Barbeau, K. A.; Allen, A. E. Carbonate-sensitive phytotransferrin 1497 controls high-affinity iron uptake in diatoms. *Nature* **2018**, 555, 534. 1498
- (89) Radivojac, P.; Clark, W. T.; et al. A large-scale evaluation of 1499 computational protein function prediction. *Nat. Methods* **2013**, *10* 1500 (3), 221–227.
- (90) Webb, E. A.; Moffett, J. W.; Waterbury, J. B. Iron Stress in 1502 Open Ocean Cyanobacteria (*SynechococcusTrichodesmium* and *Croco-* 1503 *sphaera*): Identification of the IdiA protein. *Appl. Environ. Microbiol.* 1504 **2001**, *67*, 5444–5452.
- (91) Hitzler, P.; Janowicz, K. Semantic Web. 3rd ed.; J. Chapman and 1506 Hall/CRC: 2014; Vol. 50, p 50–10–13.
- (92) Patton, E. W.; Seyed, P.; Wang, P.; Fu, L.; Dein, F. J.; Bristol, R. 1508 S.; McGuinness, D. L. SemantEco: A semantically powered modular 1509 architecture for integrating distributed environmental and ecological 1510 data. Future Generation Computer Systems 2014, 36, 430–440.
- (93) Martens, L.; Hermjakob, H.; Jones, P.; Adamski, M.; Taylor, 1512 C.; States, D.; Gevaert, K.; Vandekerckhove, J.; Apweiler, R. PRIDE: 1513 The proteomics identifications database. *Proteomics* **2005**, *5* (13), 1514 3537–3545.
- (94) Pino, L. K.; Searle, B. C.; Huang, E. L.; Noble, W. S.; 1516 Hoofnagle, A. N.; MacCoss, M. J. Calibration Using a Single-Point 1517 External Reference Material Harmonizes Quantitative Mass Spec- 1518 trometry Proteomics Data between Platforms and Laboratories. *Anal.* 1519 *Chem.* 2018, 90 (21), 13112–13117.
- (95) Kleiner, M., Normalization of metatranscriptomic and 1521 metaproteomic data for differential gene expression analyses: The 1522 importance of accounting for organism abundance. *PeerJ. Preprints* 1523 **2017**, *5*, e2846v1
- (96) Johnson, Z. I.; Zinser, E. R.; Coe, A.; McNulty, N. P.; 1525 Woodward, E. M. S.; Chisholm, S. W. Niche Partitioning Among 1526 *Prochlorococcus* Ecotypes Along Ocean-Scale Environmental Gra- 1527 dients. *Science* **2006**, 311 (5768), 1737–1740.
- (97) White, F. M. The potential cost of high-throughput proteomics. 1529 Sci. Signaling 2011, 4 (160), pe8.

Journal of Proteome Research

1531 (98) Doney, S. C. The Growing Human Footprint on Coastal and 1532 Open-Ocean Biogeochemistry. *Science* **2010**, 328 (5985), 1512–1516.