

1 Ribosomal protein phylogeography offers quantitative
2 insights into the efficacy of genome-resolved surveys of
3 microbial communities

4 Matthew S. Schechter^{1,2,*}, Florian Trigodet^{3,4}, Iva A. Veseli^{3,4}, Samuel E. Miller⁵, Matthew
5 L. Klein², Metehan Sever^{3,4}, Loïs Maignien⁶, Tom O. Delmont^{7,8}, Samuel H. Light^{2,9,*†},
6 A. Murat Eren^{3,4,5,10,11,*†}

7 ¹Committee on Microbiology, University of Chicago, Chicago, IL 60637, USA; ²Duchossois Family Institute, University
8 of Chicago, Chicago, IL 60637, USA; ³Helmholtz Institute for Functional Marine Biodiversity, 26129 Oldenburg,
9 Germany; ⁴Alfred Wegener Institute, Helmholtz Centre for Polar and Marine Research, 27570 Bremerhaven, Germany;
10 ⁵Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods
11 Hole, MA 02543, USA; ⁶University of Brest, CNRS, IFREMER, EMR 6002 BIOMEX, Unité Biologie et Écologie des
12 Écosystèmes Marins Profonds BEEP, F-29280 Plouzané, France; ⁷Génomique Métabolique du Genoscope, Institut
13 François Jacob, CEA, CNRS, University of Évry Val d'Essonne, Université Paris-Saclay, Evry, France; ⁸Research
14 Federation for the Study of Global Ocean Systems Ecology and Evolution, FR2022/Tara GOSEE, Paris, France;
15 ⁹Department of Microbiology, University of Chicago, Chicago, IL 60637, USA; ¹⁰Institute for Chemistry and Biology of
16 the Marine Environment, University of Oldenburg, 26129 Oldenburg, Germany; ¹¹Max Planck Institute for Marine
17 Microbiology, 28359 Bremen, Germany.

18 * Correspondence: meren@hifmb.de, samlight@uchicago.edu, and mschechter@uchicago.edu

19 † Co-senior authors: samlight@uchicago.edu and meren@hifmb.de

20 **Running Title**

21 Ribosomal protein phylogeography with EcoPhylo

22 **Keywords**

23 phylogeography, metagenomics, metagenome-assembled genomes, single amplified genomes

24 **Supplementary Information**

25 doi:[10.6084/m9.figshare.28200050](https://doi.org/10.6084/m9.figshare.28200050)

26 Abstract

27 The increasing availability of microbial genomes is essential to gain insights into microbial ecology
28 and evolution that can propel biotechnological and biomedical advances. Recent advances in
29 genome recovery have significantly expanded the catalogue of microbial genomes from diverse
30 habitats. However, the ability to explain how well a set of genomes account for the diversity in a
31 given environment remains challenging for individual studies or biome-specific databases. Here
32 we present EcoPhylo, a computational workflow to characterize the phylogeography of any gene
33 family through integrated analyses of genomes and metagenomes, and our application of this
34 approach to ribosomal proteins to quantify phylogeny-aware genome recovery rates across three
35 biomes. Our findings show that genome recovery rates vary widely across taxa and biomes, and
36 that single amplified genomes, metagenome-assembled genomes, and isolate genomes have
37 non-uniform yet quantifiable representation of environmental microbes. EcoPhylo reveals highly
38 resolved, reference-free, multi-domain phylogenies in conjunction with distribution patterns of
39 individual clades across environments, providing a means to assess genome recovery in
40 individual studies and benchmark biome-level genome collections.

41 Introduction

42 Establishing comprehensive genome catalogues is a fundamental objective in microbiology as
43 genomes are essential to develop insights into microbial life and to advance biotechnology and
44 biomedicine (Eren and Banfield 2024). Indeed, the rapidly increasing number of microbial
45 genomes (1) provides an evolutionary framework to resolve the branches of the Tree of Life (C.
46 T. Brown et al. 2015; Spang et al. 2015), (2) enables hypothesis generation and testing through
47 comparative genomics (Paoli et al. 2022; Al-Shayeb et al. 2022; Durrant et al. 2023; J. Chen et
48 al. 2024a), (3) offers resources to search for novel biosynthetic capabilities and natural products
49 (Paoli et al. 2022; J. Chen et al. 2024b), (4) contributes to the body of nucleotide data used to
50 train biological language models (Comman et al. 2024; Nguyen et al. 2024; Hwang et al. 2024)
51 and more, while well-structured databases aim to consolidate and give access to the outcomes
52 of genome recovery efforts (Parks et al. 2022; Schmidt et al. 2024).

53 Increasing availability of microbial genomes is a result of multiple complementary breakthroughs
54 that include (1) advances in high-throughput or targeted cultivation that enable the recovery of
55 isolate genomes (Jiang et al. 2016; Watterson et al. 2020; Cross et al. 2019), (2) the use of
56 environmental shotgun sequencing that enables the recovery of metagenome-assembled
57 genomes (MAGs) (L.-X. Chen et al. 2020), and (3) the use of microfluidics and cell sorting that
58 enables the recovery of single amplified genomes (SAGs) (Woyke, Doud, and Schulz 2017).
59 These strategies have not only been used in large-scale characterization of many of the Earth's
60 biomes (Pasolli et al. 2019; Parks et al. 2017; Pachiadaki et al. 2019; Ma et al. 2023), but also
61 have been applied to many specific questions or niche systems that span a wide range of research
62 priorities, collectively resulting in over 500,000 non-redundant bacterial and archaeal genomes
63 (Parks et al. 2022). The recovery of microbial genomes is now a relatively well-established
64 practice, yet it is not straightforward to assess (1) how taxonomic or biome-specific biases impact

65 on genome recovery efforts, and (2) the ecological or evolutionary importance of unrecovered
66 populations. As a result, individual studies that recover genomes, or efforts that curate biome-
67 specific or global genomic collections, rarely offer quantitative insights into one of the key
68 questions they aim to address: “how well do these genomes represent this environment?”.

69 Attempts to benchmark genome recovery often rely upon metagenomic read recruitment statistics
70 to quantify the fraction of reads that map to genomes with the assumption that the proportion of
71 reads recruited by a genomic collection is a proxy for the degree to which a genome collection
72 represents the genomic fragments found in a given environment. In individual studies that
73 reconstruct genomes directly from environmental metagenomes, the proportion of metagenomic
74 reads that are recruited by resulting MAGs can vary from as low as 7% in the surface ocean
75 (Delmont et al. 2018a) to as high as 80% in the human gut (Carter et al. 2023). While read
76 recruitment statistics are easy to generate and communicate, they fail to contextualize what is
77 present in the unmapped fraction and thus leave considerable ambiguity about the microbial
78 community. For instance, a large fraction of metagenomic reads not mapping to the genome
79 catalogue could belong to a single organism or multiple taxonomically diverse microbes with
80 critical ecological roles in the system. Furthermore, genome collections often systematically
81 underrepresent certain portions of the tree of life, as the rate of genome recovery differs across
82 taxa as a function of genome recovery methodology: while cultivation efforts often struggle to
83 capture slow-growing organisms (Imachi et al. 2020) or those that depend on others for survival
84 (He et al. 2015), genome-resolved metagenomics often struggle to reconstruct genomes from
85 taxa that form highly complex populations (Giovannoni 2017; Pachiadaki et al. 2019). Altogether,
86 biological and non-biological factors confound accurate interpretations of read recruitment results,
87 and the ability to measure genome recovery rates requires alternative strategies that can
88 contextualize the ecological and evolutionary relationships of organisms recovered in genome
89 collections with environmental populations accessible through metagenomics.

90 One approach to gaining insight into microbial life underrepresented in genome catalogues
91 involves the use of marker genes. *De novo* assembly, in which individual sequencing reads are
92 stitched together to recover much longer contiguous segments of DNA (contigs), is common to
93 the vast majority of genome recovery efforts. While in most cases contigs only represent
94 fragments of genomes, they still explain a much greater genomic context than unassembled reads
95 and give access to entire open reading frames, including phylogenetically informative marker
96 genes. Employing such phylogenetically informative genes assembled from metagenomes in
97 conjunction with metagenomic read recruitment enables fine-grained analyses of phylogeny and
98 biogeography of individual taxa, as demonstrated by previous studies that used the *rpoC1* gene
99 to characterize the phylogeography of marine bacteria (Kent et al. 2019; Ustick, Larkin, and
100 Martiny 2023) or RNA and DNA polymerases to identify and guide the genomic recovery of major
101 viral clades (Weinheimer and Aylward 2020; Gaia et al. 2023).

102 Among all phylogenetically informative genes, ribosomal proteins represent a special class as
103 they (1) occur as a single-copy gene in genomes across the tree of life, (2) are consistently
104 assembled even for complex or relatively rare populations in metagenomes due to their relatively
105 short length, and (3) contain enough phylogenetic information to delineate distinct branches of life
106 at relatively high levels of resolution (Olm et al. 2020). Recognizing their utility, many studies have
107 leveraged individual ribosomal proteins to analyze community composition (Wu and Eisen 2008;
108 Crits-Christoph et al. 2022), integrating ribosomal protein phylogenies with metagenomic read
109 recruitment to track individual clades of microbes (Hug et al. 2013; Emerson et al. 2016; Hamilton
110 et al. 2016; Diamond et al. 2019; Matheus Carnevali et al. 2021). Ribosomal proteins are thus
111 ideally suited gene markers for tracking microbial populations underrepresented within genome
112 collections.

113 Here we present EcoPhylo, a workflow to simultaneously visualize the phylogenetic relationships
114 and biogeographical distribution patterns of sequences that match any given gene family from

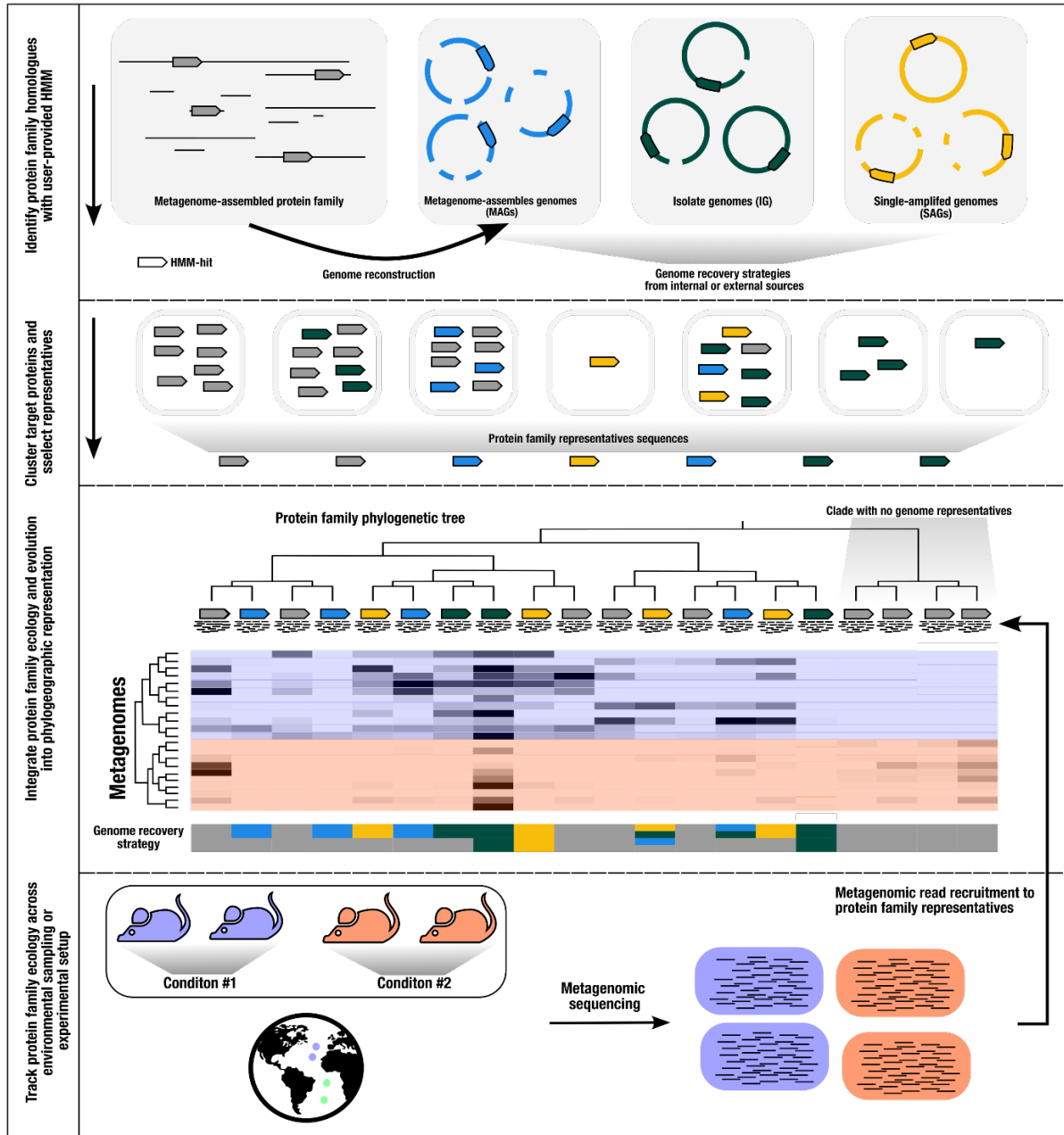
115 genomes and metagenomes, and demonstrate its application to the phylogeography of ribosomal
116 proteins for quantification of genome recovery rates across biomes. Our results show that bringing
117 together multi-domain ribosomal protein phylogenies with distribution patterns of individual clades
118 across environments in a single interface offers a valuable data analysis and visualization strategy
119 to benchmark genome recovery efforts scaling from individual projects to global surveys of large
120 genome collections and metagenomes.

121 Results

122 EcoPhylo enables integrated surveys of gene family 123 phylogeography

124 EcoPhylo implements a computational workflow to integrate the phylogeny and biogeography of
125 any given gene family and enables its users to track the distribution patterns and evolutionary
126 relationships between homologous genes across environments and/or experimental conditions
127 (Figure 1, also see Materials and Methods).

128 When applied to phylogenetically tractable single-copy core genes, such as ribosomal proteins,
129 in tandem with metagenomes and a genome collection, EcoPhylo identifies populations
130 assembled in metagenomes but absent in the genomic collections (and *vice versa*), highlighting
131 the ecological and evolutionary relevance of organisms detected through metagenomic
132 assemblies but lacking genomic representation (Figure 1). This allows for the quantification of
133 genome recovery rates of different methods (e.g., isolate genomes, MAGs, SAGs) across taxa
134 and provides a means to investigate phylogenetic and ecological features of organisms without
135 genomic representation. Importantly, the unbroken link between genes and contigs enables
136 downstream targeted binning efforts when necessary.



137

138 **Figure 1: Schematic of the EcoPhylo workflow applied to a single gene family.** The proposed workflow
 139 integrates biogeography from metagenomic read recruitment and protein phylogenetics to display the
 140 phylogeographical distribution of closely related lineages. When including genome sources, the workflow highlights
 141 which genome recovery strategies are more effective for sampling specific taxa. Although this manuscript focuses on
 142 ribosomal proteins, the proposed workflow is generalizable to any gene family.

143 Using ribosomal proteins to *de novo* characterize the phylogenetic makeup of microbiomes and
 144 benchmark genome recovery rates has numerous advantages. However, these advantages also
 145 pose noteworthy challenges. Ribosomal proteins are short protein sequences (~300 amino acids),

146 which substantially limits their ability to resolve deep phylogenetic branching patterns.
147 Furthermore, their evolution is subject to strong purifying selection, as a result, the average
148 nucleotide identity (ANI) threshold often used to define 'species' boundaries between whole
149 genomes is 95% (Jain et al. 2018) increases to 99% for ribosomal protein sequences (Olm et al.
150 2020). Therefore, ribosomal proteins are more vulnerable than other genes to non-specific read-
151 recruitment from closely related proteins within metagenomes. To identify criteria for reliably
152 resolving taxa, we started our investigation by developing a series of benchmarks to optimize the
153 use of ribosomal proteins in EcoPhylo with appropriate parameters to maximize the ecological
154 and evolutionary signal they can offer while minimizing non-specific read recruitment. These
155 benchmarks, which are detailed in the Supplementary information (1) inspected hidden Markov
156 model (HMM) alignment coverage thresholds to accurately detect ribosomal proteins in genomes
157 and metagenomes; (2) examined the copy number distribution of ribosomal protein HMMs across
158 archaeal and bacterial genomes to only consider single-copy candidates; and (3) explored
159 nucleotide similarity thresholds to cluster ribosomal gene sequences to maximize the taxonomic
160 resolution of representative sequences while maintaining sufficient nucleotide distance between
161 distinct representative sequences to reduce non-specific read recruitment from metagenomes
162 (Supplementary information).

163 Based on these considerations, we implemented routines and adjusted default EcoPhylo
164 parameters to (1) use a minimum of 80% model coverage for ribosomal protein HMMs for a match;
165 (2) filter for complete open reading frame sequences to remove assembly artifacts; and (3) cluster
166 HMM hits with target coverage to ensure grouping of extended open reading frames and leverage
167 97% nucleotide similarity as the most appropriate clustering threshold to minimize non-specific
168 read recruitment (Supplementary information). We also compared broad ecological insights
169 recovered from EcoPhylo to state-of-the-art taxonomic profiling tools, confirming that this
170 framework offered qualitatively comparable results (Supplementary information). Altogether,

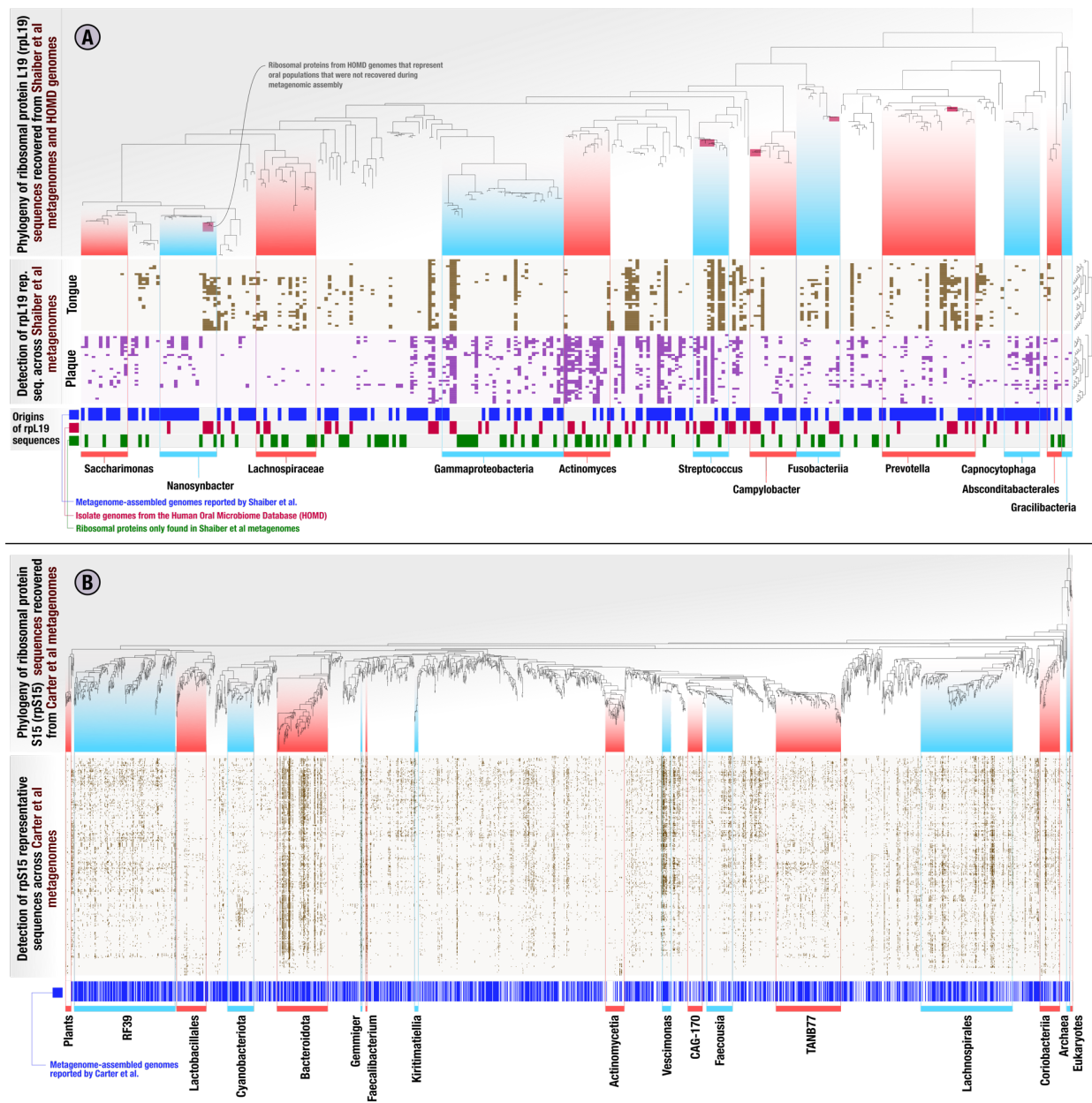
171 these evaluation and optimization steps yielded EcoPhylo default parameters to obtain
172 representative ribosomal protein sequences that are suitable for investigations of the phylogeny,
173 biogeography, and genome recovery of populations they describe.

174 Ribosomal proteins quantify and contextualize genome recovery 175 rates from metagenomes

176 Thanks to its diverse physiological properties that promote a variety of chemical gradients and
177 surfaces (Bowen et al. 2018), the human oral cavity is home to diverse communities of microbes
178 (Dewhirst et al. 2010). The human oral microbiome is a relatively well-characterized environment
179 with a wealth of isolate genomes accessible through the Human Oral Microbiome Database
180 (HOMD) (Escapa et al. 2018; T. Chen et al. 2010), and numerous genome-resolved
181 metagenomics surveys that have captured representative genomes of microbial clades that have
182 largely eluded cultivation efforts. Using EcoPhylo we first focused on a genome-resolved
183 metagenomics survey which reconstructed multiple high-quality MAGs from tongue and plaque
184 samples from the human oral cavity (Shaiber et al. 2020). While Shaiber et al. (2020) reported
185 numerous genomes for elusive taxa, such as *Saccharimonadia* (TM7), *Absconditabacteria* (SR1),
186 and *Gracilibacteria* (GN02), the genome-resolved metagenomic workflow failed to reconstruct
187 MAGs that resolved to some of the best-represented organisms in culture collections from the
188 oral cavity, such as members of the genus *Streptococcus*, (Escapa et al. 2018), which was
189 represented by only two MAGs in Shaiber et al. (2020). This discrepancy compelled us to combine
190 isolate genomes from the HOMD together with metagenomes and MAGs from Shaiber et al.
191 (2020), to investigate whether EcoPhylo could reveal the differential recovery of genomes through
192 distinct recovery approaches.

193 We started our analysis by combining 790 non-redundant MAGs and 14 metagenomic co-
194 assemblies of tongue and plaque metagenomes reported by Shaiber et al (2020) with 8,615
195 isolate genomes we obtained from the HOMD (Supplementary Table 1). To characterize these
196 data, we elected to use EcoPhylo with *rpL19* HMM, since it was the most frequent ribosomal
197 protein with an average length of 393 nucleotides across all genomes in our collection, occurring
198 in 98.59% of the HOMD genomes and 81.81% of the Shaiber et al. MAGs (Supplementary
199 information). To assess the generalizability of observations made from *rpL19*, we also ran
200 EcoPhylo on the same dataset with *rpS15* and *rpS2*, with the average length of 275 and 781
201 nucleotides, respectively (Supplementary Table 2, Supplementary information).

202 The EcoPhylo analysis of the *rpL19* genes found in the genomes and metagenomic assemblies
203 resulted in a phylogenetic tree with 277 non-redundant bacterial representative sequences
204 (Figure 2A, Supplementary Table 3). Hierarchical clustering of metagenomes based on the
205 detection patterns of these *rpL19* sequences organized metagenomes into tongue and plaque
206 sampling sites *de novo* (Figure 2A, Supplementary Figure 1), demonstrating that a single
207 ribosomal gene family is able to capture the known ecological differences between these habitats.
208 Many closely related *rpL19* genes that resolved to prevalent oral taxa, such as *Prevotella* and
209 *Streptococcus*, showed within-genus differences in site specificity, a previously observed
210 phenomenon (Eren et al. 2014) that is attributed to divergent accessory genomes (Mark Welch,
211 Dewhirst, and Borisy 2019; Utter et al. 2020). Multiple ribosomal protein representative sequences
212 recruited reads from tongue as well as plaque metagenomes, also matching prior observations of
213 cosmopolitan taxa (Figure 2A, Supplementary information). Overall, the ecological insights
214 revealed by *rpL19* recapitulated known ecology of oral microbes (Mark Welch, Dewhirst, and
215 Borisy 2019) and provided a framework to assess genome recovery rates.



216

217 **Figure 2: Ribosomal protein phylogeny and detection patterns across metagenomes from the human oral**

218 **cavity and gut microbiomes.** In the heatmaps in both panels, each column represents a ribosomal protein

219 representative sequence, each row represents a metagenome, and each data point indicates whether a given ribosomal

220 protein was detected in a given metagenome. The columns of heatmaps are ordered by a tree which represents a

221 phylogenetic analysis of all ribosomal protein representative sequences, and the rows are ordered by a hierarchical

222 clustering dendrogram that is calculated based on the ribosomal protein detection patterns across metagenomes. The

223 panel **(A)** represents the EcoPhylo analysis of *rpL19* sequences across Shaiber et al. (2020) metagenome-assembled

224 genomes (MAGs), Shaiber et al. (2020) oral metagenomes, and includes three additional rows

225 that indicate the origin of a given ribosomal protein, whether it is a metagenome-assembled genome (MAG, blue),
226 HOMD isolate genome (red), or only recovered from metagenomic assemblies with no representation in genomes
227 (green). Smaller red boxes in the phylogenetic tree mark microbial clades that were absent in the collection of MAGs
228 and assemblies reported by Shaiber et al. (2020), but detected in Shaiber et al. (2020) metagenomes solely due to the
229 inclusion of HOMD isolate genomes. The panel **(B)** represents the EcoPhylo analysis of *rpS15* sequences across the
230 Carter et al. (2023) metagenome-assembled genomes (MAGs) and Carter et al. (2023) gut metagenomes from a Hadza
231 tribe, and includes an additional row that indicates whether a MAG was reported for a given ribosomal protein (blue).

232 EcoPhylo tracks the origins of each sequence in each sequence cluster. Some *rpL19* clusters,
233 representatives of which are shown in the phylogenetic tree in Figure 2A, were composed of
234 sequences found only in metagenomic assemblies and not in MAGs or isolate genomes,
235 highlighting clades present in the environment but not in genome collections. Other *rpL19* clusters
236 only contained sequences represented in HOMD isolate genomes; despite their consistent
237 detection in oral samples through metagenomic read recruitment, they were absent in
238 metagenomic assemblies or MAGs, highlighting clades that are less accessible to short-read
239 metagenomic assembly approaches (Figure 2). To calculate genome recovery rates for any given
240 taxon, we divided the number of sequence clusters that contained a sequence from a given
241 genome recovery method by the total number of representative sequences EcoPhylo reported for
242 that taxon (Materials and Methods). This analysis revealed that 60.3% of the bacterial populations
243 defined by *rpL19* gene clusters that were detected in metagenomic reads also appeared in MAGs.
244 In other words, the overall bacterial MAG recovery rate in the study by Shaiber et al. (2020) was
245 60.3% (*rpS15*: 62.8%, *rpS2*: 53.2%) (Figure 2A, Supplementary Table 3, Supplementary Table
246 4). However, this rate of recovery was not uniform across individual taxa. EcoPhylo revealed
247 higher MAG recovery rates for taxa such as *Saccharimonas* at 69.2% (*rpS15*: N/A, *rpS2*: 63.6%),
248 and *Prevotella* at 76.9% (*rpS15*: 82.6%, *rpS2*: 84%). In contrast, the MAG recovery was lower for
249 populations in other clades, including *Gammaproteobacteria* and *Fusobacteriia*, with MAG
250 recovery rates of 47.1% (*rpS15*: 58.1%, *rpS2*: 44.8%) and 41.7% (*rpS15*: 46.2%, *rpS2*: 31.2%),

251 respectively (Figure 2A, Supplementary Table 3, Supplementary Table 4). The MAG recovery
252 rate was particularly low for *Streptococcus* at 30% (*rpS15*: 15.4%, *rpS2*: 10%), consistent with
253 the presence of only two MAGs in Shaiber et al. (2020). However, the MAG recovery rate for
254 *Actinomyces* was also very low at 23.1% (*rpS15*: 36.4%, *rpS2*: 13.3%) despite the
255 characterization of nine *Actinomyces* MAGs by Shaiber et al. (2020) reveals a large number of
256 distinct *Actinomyces* populations missed by MAGs even though they were present in the
257 assemblies (Figure 2A, Supplementary Table 3, Supplementary Table 4). Overall, this analysis
258 not only confirmed that MAG recovery rates are not uniform across microbial clades, but also
259 showed that quantification of these rates is possible and may yield unexpected insights into the
260 extent of diversity that is not represented in the final set of MAGs for some clades.

261 The inclusion of genomes from the HOMD increased the number of *rpL19* sequence clusters that
262 contained genomes in this dataset, i.e., the total genome recovery rate, from 60.3% to 73.3%
263 (*rpS15*: 74.8%, *rpS2*: 81.3%), and led to the representation of 35 additional microbial clades for
264 which the metagenomic sequencing and analysis workflow implemented in Shaiber et al. (2020)
265 did not assemble. As with MAGs, the improved detection of taxa among HOMD genomes was
266 not uniform across clades (Figure 2A). For example, HOMD genomes offered genomic context
267 for five additional *Streptococcus* populations, increasing the genome recovery rate from 30% with
268 MAGs only, up to 80% when including the HOMD collection. When taking into account both MAGs
269 and isolate genomes, the overall genome recovery rate of Shaiber et al. (2020) from a human
270 oral microbiome dataset determined by EcoPhylo was 73.3%, showing that ribosomal protein
271 phylogeography is an effective means to quantify genome recovery statistics for individual
272 studies. Conversely, EcoPhylo results showed that 26.7% of the individual clades that could be
273 detected through the presence of *rpL19* sequences in assemblies of Shaiber et al. (2020)
274 metagenomes lacked genomic representation in both Shaiber et al. (2020) MAGs and HOMD
275 isolates (Figure 2A). Clades that were solely detected through their assembled yet not binned

276 ribosomal proteins increased the detection of populations of *Lachnospiraceae*, *Actinomyces*,
277 *Gammaproteobacteria*, and *Patescibacteria* (Figure 2A). As EcoPhylo clusters ribosomal proteins
278 at 97% nucleotide similarity, a conservative threshold that underestimates biodiversity by often
279 grouping genomes with gANI below 95% (Olm et al. 2020).

280 Next, we applied EcoPhylo to another genome-resolved metagenomics study that recently
281 characterized the gut microbiome of a Hadza hunter-gatherer tribe with a deep sequencing effort
282 by Carter et al. (2023), in which the authors reported nearly 50,000 redundant bacterial and
283 archaeal MAGs from 338 metagenomes with an average of 76 million paired-end reads
284 (Supplementary Table 1). EcoPhylo analysis of this dataset with *rpS15* with an average length of
285 276 nucleotides, along with *rpS16* and *rpL19*, with the average length of 297 and 370 nucleotides
286 respectively (Supplementary Figure 2), revealed a relatively high bacterial MAG recovery rate of
287 67.7% (*rpS16*: 72.8%, *rpL19*: 69.5%) (Figure 2B, Supplementary Table 2). While there were some
288 clades, such as *Actinomycetia*, for which the genome recovery rate was as low as 31.9% (*rpS16*:
289 32.4%, *rpL19*: 33.3%), the high MAG recovery rate was generally uniform across all major taxa
290 (Supplementary Table 5, Supplementary Table 6, Supplementary Information).

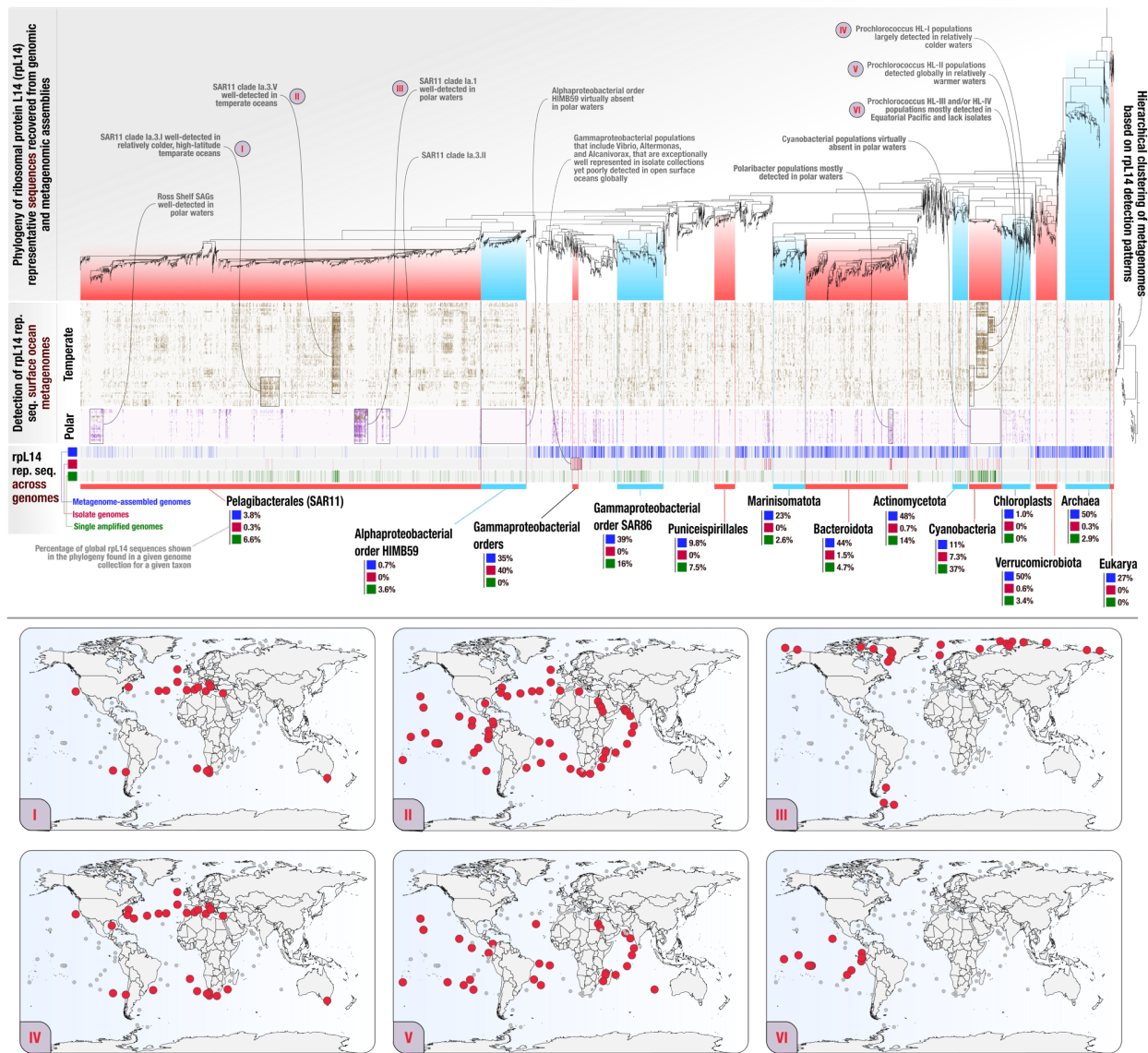
291 Through these analyses, we are able to demonstrate that the MAGs obtained by Carter et al.
292 (2023) more comprehensively represents the populations captured by their metagenomic
293 assemblies of the human gut compared to the MAGs obtained by Shaiber et al. (2020) given their
294 metagenomic assemblies of the oral cavity (Figure 2B, Supplementary Table 5, Supplementary
295 Table 6, Supplementary Information). The ability to make such a statement highlights the utility of
296 EcoPhylo at providing quantitative insights into the efficacy of genome-resolved surveys
297 independent of biomes while offering a phylogenetic and biogeographical context for the
298 populations that were detected in the assemblies.

299 Overall, EcoPhylo results from the human oral cavity and human gut ecosystems show that our
300 workflow can scale to large metagenomic surveys, combine genomes from multiple sources to
301 compare distinct recovery strategies at the level of individual phylogenetic clades, and
302 recapitulate known ecological patterns.

303 **Genome collections represent a small fraction of microbial diversity** 304 **in the global surface ocean microbiome**

305 Marine systems support fundamental biogeochemical cycles that maintain the Earth's habitability,
306 and comprehensively documenting the genomes of marine microbes that are intimately
307 connected to these processes has been one of the key aims of microbiology. In addition to
308 decades of cultivation efforts, recent years witnessed a rapid expansion of marine microbial
309 genome catalogues for bacteria and archaea with new MAGs (Delmont et al. 2018b; Tully,
310 Graham, and Heidelberg 2018; Paoli et al. 2022) and SAGs (Pachiadaki et al. 2019; Martínez-
311 Pérez et al. 2022). Studies that recover genomes from marine systems recognize that the extent
312 to which these collections represent marine environmental populations is limited (Delmont et al.
313 2019a; Paoli et al. 2022). Yet, quantifying the extent of representation at the level of individual
314 environmental clades across genome collections is a challenge. Having established the utility of
315 EcoPhylo to elicit quantitative answers to such questions, we next surveyed a state-of-the-art
316 globally distributed collection of microbial genomes from marine systems (Paoli et al. 2022) in the
317 context of metagenomes generated by the Tara Oceans Project (Salazar et al. 2019; Sunagawa
318 et al. 2015), the Hawaii Ocean Time-series (HOT) (Biller et al. 2018), the Bermuda Atlantic Time-
319 series (BATS) (Biller et al. 2018), BioGEOTRACES expeditions (Biller et al. 2018), and the
320 Malaspina Project (Sánchez et al. 2024) to simultaneously compare genome recovery rates of
321 MAGs, SAGs, and isolate genomes. Of all 1,038 metagenomes, we focused on those that were
322 collected from up to 30m depth and had a size fraction of 0.22 μ m to 3 μ m (Supplementary Figure

323 3, Supplementary Table 1), which left us with a total of 237 metagenomes containing a total of
324 18,832,767,852 short reads (79,463,155 reads per metagenome on average). Our collection of
325 genomes included 7,282 MAGs, 1,474 SAGs, and 1,723 isolate genomes from The Ocean
326 Microbiomics Database (subsetting from samples of < 30m depth when possible) (Paoli et al.
327 2022). We expanded this collection with an additional 52 isolate genomes that historically have
328 low MAG recovery rates, such as *Pelagibacterales* (SAR11) and *Cyanobacteriota*, and a
329 collection of 41 SAGs obtained from below the Ross Ice Shelf to improve detection of cold-
330 adapted clades (Martínez-Pérez et al. 2022) (Materials and Methods), yielding a total of 10,479
331 genomes (Supplementary Table 1). For characterization of these data by EcoPhylo we primarily
332 used the ribosomal gene *rpL14* with an average length of 363 nucleotides, which we detected in
333 82% of the final list of genomes (Supplementary information), but we also conducted additional
334 analyses using the ribosomal genes *rpS8* and *rpS11*, with an average length of 398 and 415
335 nucleotides respectively, to confirm our key observations (Supplementary Figure 1,
336 Supplementary Table 2).



337

338 **Figure 3: Ribosomal protein L14 phylogeny and detection patterns across metagenomes from the global**

339 **surface ocean (depth < 30 m; size fraction: 0.22 to 1.6 μm , 0.22 to 3 μm).** In the heatmap of panel (A), each column

340 represents a ribosomal protein representative sequence, each row represents a metagenome, and each data point

341 indicates whether a given ribosomal protein was detected in a given metagenome. The heatmap columns are ordered

342 by a tree which represents a phylogenetic analysis of all ribosomal protein representative sequences, and the rows are

343 ordered by a hierarchical clustering dendrogram that is calculated based on ribosomal protein detection patterns across

344 metagenomes. Metagenomes are colored by temperate (gold) or polar (purple) biomes. Each leaf of the phylogenetic

345 tree is decorated below the heatmap with metadata denoting the origin of the RP: metagenome-assembled genome

346 (MAG) (blue), isolate genomes (red), and single amplified genomes (green). In panel (B), each map corresponds to

347 phylogeographical patterns highlighted in panel (A). Colored sampling points correspond to the boxed phylogeographic
348 signals in panel (A).

349 EcoPhylo analysis of *rpL14* genes across 236 global surface ocean metagenomes characterized
350 8,075 bacterial, 370 archaeal, and 33 eukaryotic clades and computed their distribution patterns
351 across environments (Supplementary Table 7). Hierarchical clustering of metagenomes based on
352 *rpL14* detection patterns split samples into two major groups, whereby one of the groups
353 represented samples collected from polar regions and the other represented samples collected
354 from temperate oceans (Figure 3, Supplementary information); a result that is in line with previous
355 observations that documented water temperature as a major driver of microbial diversity in the
356 surface ocean (Sul et al. 2013; Sunagawa et al. 2015). Notably, temperate and polar water
357 samples did not partition when we included metagenomes with lower sequencing depths in our
358 analysis. This was likely caused by increasing noise in detection patterns of various prevalent
359 populations, which compelled us to only consider metagenomes with 50 million or more paired-
360 end reads for our downstream analyses (Supplementary information), which left us with a total of
361 100 metagenomes (Supplementary Table 1, Supplementary information). Overall, EcoPhylo
362 captured (1) differential distribution patterns among closely related taxa as a function of
363 temperature and latitude, a form of phylogenetic overdispersion likely due to greater competitive
364 exclusion among closely related organisms in the same ecological niche, and (2) showed that the
365 majority of taxa contained both warm- and cold-adapted clades that exclusively occurred either in
366 polar or temperate waters, an expected observation since marine thermal adaptation is not
367 correlated with phylogenetic signal (Thomas et al. 2012) and is likely acquired through
368 independent processes within each major clade (Figure 3). Furthermore, the *rpL14*
369 phylogeography captured well-understood biogeographical patterns of prevalent pelagic taxa
370 (Figure 3), in agreement with previous studies that showed the dominance of SAR11 subclade
371 Ia.3.V in temperate waters (Delmont et al. 2019b) and the exclusivity of cold-adapted SAR11
372 clades Ia.1 and Ia.3.II to polar regions (M. V. Brown et al. 2012; Delmont et al. 2019a). It also

373 corroborated the global distribution of *Prochlorococcus* HL-II in temperate waters (Johnson et al.
374 2006; Biller et al. 2015; UstICK, Larkin, and Martiny 2023) and the contrasting distribution of this
375 group with *Prochlorococcus* HL-III and *Prochlorococcus* HL-IV, which are mainly found in the
376 Equatorial Pacific (Rusch et al. 2010; Huang et al. 2012; Malmstrom et al. 2013; Kent et al. 2016),
377 as well as *Prochlorococcus* HL-I, which is confined to higher latitudes (Johnson et al. 2006; Biller
378 et al. 2015; Delmont and Eren 2018). The concordance of our results from EcoPhylo with known
379 ecological patterns in marine microbiology underscores the reliability of ribosomal protein
380 phylogeography in characterizing the interplay between microbial ecology and evolution in global
381 surface ocean microbiome (Figure 3, Supplementary information).

382 Using these data, we first compared the overlap between surface ocean microbial populations in
383 the environment and publicly available MAGs generated from this biome by calculating genome
384 recovery rates. The MAG recovery rate for Archaea was relatively high at 49.5% (*rpS8*: N/A,
385 *rpS11*: 50%), however, the MAG recovery rate for Bacteria was only 19.9% (*rpS8*: 22.7%, *rpS11*:
386 22.2%). In contrast to the MAG recovery rates we observed in individual studies from the human
387 oral and gut microbiome (60.3% and 67.7%, respectively), this much lower recovery rate from
388 multiple sequencing projects reflects the relatively poor efficiency and the contemporary
389 challenges of reconstructing genomes from metagenomes in the ocean biome (Figure 3,
390 Supplementary Figure 1, Supplementary Table 7, Supplementary Table 8). Some phyla had
391 relatively high MAG recovery rates, such as 48.2% for *Actinomycetota* (*rpS8*: 50.8%, *rpS11*:
392 48.9%), 43.9% for *Bacteroidota* (*rpS8*: 47.3%, *rpS11*: 44.3%), and 49.7% for *Verrucomicrobiota*
393 (*rpS8*: 51.2%, *rpS11*: 44.4%). MAG recovery rates were much lower for other clades, including
394 those containing some of the best-studied autotrophs and heterotrophs of open surface oceans,
395 for example, the MAG recovery rate was only 11% for *Cyanobacteriota* (*rpS8*: 11.7%, *rpS11*:
396 9.93%). Many clades of *Alphaproteobacteria* had some of the lowest MAG recovery rates,
397 including 12.7% (*rpS8*: 9.21%, *rpS11*: 7.2%) for the uncharacterized order HIMB59

398 (Supplementary Table 7, Supplementary Table 8). Poor MAG recovery rate was also true for the
399 order *Pelagibacterales*, which remained at 3.76% (*rpS8*: 3.98%, *rpS11*: 4.25%) (Supplementary
400 Table 7, Supplementary Table 8). Compared to taxonomic classification of shotgun metagenomic
401 reads or sequencing of 16S rRNA gene amplicons, prior studies observed much lower relative
402 abundance estimates for populations resolving to *Cyanobacteriota* and *Pelagibacterales* based
403 on MAGs (Pachiadaki et al. 2019; Chang et al. 2024). By elucidating clade-specific discrepancies
404 between different methods of genome recovery, EcoPhylo offers a context for the extent of
405 missing MAGs in prior surveys, which likely is a by product of fragmented metagenomic
406 assemblies due to co-occurring closely related populations with high genomic diversity (L.-X.
407 Chen et al. 2020).

408 *De novo* characterization of *rpL14* sequences with EcoPhylo uncovers the vast diversity within
409 *Pelagibacterales* compared to the other clades ([Figure 3](#)). Strikingly, even with the conservative
410 profiling of EcoPhylo that will occasionally pull together ribosomal proteins that belong to genomes
411 from multiple 95% gANI clusters, *Pelagibacterales* made up 41.54% of the non-redundant *rpL14*
412 sequence clusters shown in [Figure 3](#), revealing yet another representation of its immense
413 phylogenetic diversity (Morris et al. 2002; M. V. Brown et al. 2012; Pachiadaki et al. 2019). While
414 both *Cyanobacteriota* and *Pelagibacterales* suffer from similar rates of poor representation in
415 MAG collections, the missing genomes for environmental populations of *Pelagibacterales*
416 resolved to ~12 times more *rpL14* sequence clusters, which unveils the enormous
417 uncharacterized genomic diversity within this order of many clades that show distinct
418 biogeographical patterns ([Figure 3](#)), and highlights the importance of ongoing cultivation efforts
419 to improve its genomic representation (Freel et al. 2024).

420 Different genome recovery methods come with different clade- 421 specific biases

422 Finally, we explored the contribution of isolate genomes and SAGs to the genomic representation
423 of surface ocean microbial populations. Isolate genomes had low phylogenetic breadth across
424 the Ribosomal L14 phylogeny and only sampled a few closely related populations, indicating the
425 repeated isolation of similar microbes. In fact, at the phylum level, isolate genomes only effectively
426 sampled *Cyanobacteriota* at a recovery rate of 7.33% (*rpS8*: 10.7%, *rpS11*: 7.80%) despite the
427 fact that we supplemented this clade with extra isolate genomes for this analysis (Supplementary
428 Table 7, Supplementary Table 8). Interestingly, a few closely related orders within class
429 *Gammaproteobacteria* were exceptionally well-covered by bacterial organisms in culture, where
430 40% of the ribosomal proteins matched to an isolate genome (Figure 3). These sister clades
431 represented a relatively small fraction of the overall phylogenetic diversity and were poorly
432 detected across the global surface ocean metagenomes, however, they collectively contained
433 many intensely studied marine model bacterial genera, such as *Vibrio* (Kauffman et al. 2018;
434 Baker-Austin et al. 2017; van Kessel and Camilli 2024; Septer and Visick 2024), *Alteromonas*
435 (Pedler, Aluwihare, and Azam 2014; Manck et al. 2022; Henríquez-Castillo et al. 2022; Z. Lu et
436 al. 2024; Halloran et al. 2025), and *Alcanivorax* (Sabirova et al. 2008; Naether et al. 2013; Manoj
437 Prasad et al. 2019; M. Prasad et al. 2023). The clades with some of the highest SAG recovery
438 rates included the order *SAR86* at 16.3% (*rpS8*: 20.1%, *rpS11*: 17.8%) and the phylum
439 *Actinomycetota* at 14.4% (*rpS8*: 16.2%, *rpS11*: 14.5%) (Figure 3, Supplementary Table 7,
440 Supplementary Table 8). Furthermore, SAGs augmented the recovery of genomes from
441 prevalent, taxonomically diverse populations with low MAG recovery rates including (1)
442 *Cyanobacteriota* with a three-fold increase compared to MAGs at 37.0% (*rpS8*: 47.2%, *rpS11*:
443 41.5%), (2) *SAR86* at 16.3% (*rpS8*: 20.1%, *rpS11*: 17.8%), and (3) *Alphaproteobacteria*, including
444 *HIMB59* at 5.06% (*rpS8*: 5.61%, *rpS11*: 5.29%) and *SAR11* at 6.35% (*rpS8*: 7.56%, *rpS11*:

445 7.82%) (Figure 3, Supplementary Table 7, Supplementary Table 8). Specifically, SAGs were able
446 to effectively sample the warm-adapted SAR11 clade 1.a.3V as well as an uncharacterized cold-
447 adapted clade of SAR11 likely due to SAG sampling sites that covered both temperate
448 (Pachiadaki et al. 2019) and polar (Martínez-Pérez et al. 2022) oceans. Considering that SAR11
449 has been estimated to be 25% of all plankton (Giovannoni 2017) and 20-40% of cells counts in
450 the surface ocean (Schattenhofer et al. 2009), SAG methodology, which separates individual
451 bacterial cells from the environment, appears to be optimal for recovering this taxon and avoids
452 the pitfalls of fragmented metagenomic assembly caused by microbiomes with closely related
453 populations (Hosokawa and Nishikawa 2024). SAGs had greater breadth than MAGs and isolate
454 genomes across the EcoPhylo phylogenies of Ribosomal L14, Ribosomal S11, and Ribosomal
455 S8 (Figure 3, Supplementary Figure 1), despite being sampled from only 9 surface ocean
456 sampling sites (Martínez-Pérez et al. 2022; Pachiadaki et al. 2019) compared to 237
457 metagenomes encompassing higher environmental diversity in the global surface ocean.
458 Furthermore, SAGs only represented 6.91% of Bacteria and 2.97% of Archaea populations
459 detected across the dataset while MAGs represented 19.9% of Bacteria and 49.5% of Archaea
460 populations detected across the dataset (Supplementary Table 7, Supplementary Table 8). These
461 results are in line with prior observations that showed pelagic SAGs represent notably more
462 taxonomic richness when compared to MAGs (Pachiadaki et al. 2019).

463 Similar to the human oral microbiome, we found an uneven phylogenetic distribution of genome
464 recovery rates among genome acquisition strategies in the global surface ocean microbiome.
465 MAGs systematically undersampled globally prevalent clades of *Alphaproteobacteria*, such as
466 SAR11. In contrast SAGs from only a few surface ocean sampling sites (n=9), substantially
467 improved their recovery rates from these clades (Figure 3), indicating that while sequencing and
468 assembly of single-cell genomes often lead to severely incomplete genomes due to amplification
469 biases (Stepanauskas et al. 2017), SAGs show great potential for unbiased genome recovery.

470 The phylogeography of ribosomal proteins adds further evidence that a combination of genome-
471 resolved metagenomics, single amplified genomics, and innovations in microbial isolation
472 strategies are needed to further increase genomic representation of diverse taxa in the global
473 surface ocean.

474 Discussion

475 Our work illuminates the efficiencies of current genome recovery methods and their ability to
476 sample genomes from various microbiomes. By leveraging phylogenetically informative marker
477 genes detected in metagenomic assemblies, such as ribosomal proteins, that are absent from
478 final genome collections, EcoPhylo provides a robust framework for benchmarking genome
479 recovery rates across multiple genome acquisition methods and contextualizing the ecological
480 and evolutionary of genome collections with naturally occurring microbial populations. Our study
481 examined three microbiome projects that used multiple genome recovery strategies (MAGs,
482 SAGs, and isolate genomes) to survey the human oral cavity, global surface ocean, and human
483 gut. Overall, we found that the EcoPhylo workflow can quantitatively measure genome recovery
484 rates and analyze heterogeneous genome collections to assess the efficacy of distinct recovery
485 methods at the level of individual phylogenetic clades. We observed that deep metagenomic
486 sequencing of the human gut microbiome yielded the highest genome recovery rate across these
487 three biomes analyzed. Additionally, we identified that a state-of-the-art genome collection from
488 marine environments represents a small fraction of the total diversity in the open surface ocean
489 through the lens of ribosomal proteins found in assembled metagenomes. By generating insights
490 into multi-domain ribosomal protein phylogeography, EcoPhylo provides a valuable interactive
491 data visualization strategy to evaluate the underlying microbial ecology of metagenomic
492 sequencing projects.

493 The *de novo* profiling of ribosomal proteins in metagenomic assemblies resembles reference-
494 based taxonomic profiling of metagenomic short reads to predict relative abundances of taxa, an
495 idea that is implemented in multiple tools that use marker genes, such as Kraken (Wood and
496 Salzberg 2014), MIDAS (Nayfach et al. 2016), Bracken (J. Lu et al. 2017), mOTUs (Ruscheweyh
497 et al. 2022), and MetaPhlAn (Manghi et al. 2023), or processed conserved marker gene windows,
498 such as SingleM (Woodcroft et al. 2024). As these tools typically report distinct taxa and their
499 relative abundances, they indeed can help assess genome recovery efforts through direct
500 comparisons of taxon names they identify to the taxonomy of recovered genomes. However, the
501 requirement of a database of reference genomes and/or marker genes, and the absence of a
502 direct link between the genes in assemblies and taxon names reported in tables limit applications
503 with additional downstream opportunities such as targeted genome recovery. In contrast, the
504 flexibility of surveying any marker gene, including ribosomal proteins, across user-provided
505 metagenomic assemblies *de novo* offers an alternative approach that directly connects genes of
506 unrecovered taxa to assemblies and estimates the number of populations detected in
507 metagenomes regardless of their phylogenetic novelty across diverse samples and conditions.

508 While the phylogeography of ribosomal proteins offers valuable insights into genome recovery,
509 these genes have notable limitations. Rates of evolution as well as the likelihood to be recovered
510 through metagenomic assembly will differ across ribosomal protein families, complicating direct
511 quantitative comparisons between different ribosomal proteins and in some cases will require
512 surveying multiple ribosomal proteins to ensure the generalizability of observations from a single
513 ribosomal protein. Additionally, individual ribosomal protein trees will have less phylogenetic
514 power compared to concatenated ribosomal protein trees or longer marker genes. Although this
515 may lead to suboptimal organism phylogenetics, the efficient organization of ribosomal proteins
516 yields informative insights into the diversity of clades within a sample. Furthermore, when working
517 with incomplete genomes, such as MAGs or SAGs, a single ribosomal gene family will rarely be

518 detected across the entire genome collection and thus only a subset of genomes will be
519 contextualized per protein. Yet the inherent trade-offs of using incomplete genomes ($x \geq 50\%$ and
520 less than 10% contamination) highlight ongoing challenges in genome recovery, as stricter
521 completeness thresholds would further reduce the number of genomes available for analysis.

522 The modular design and customizable parameters of EcoPhylo allows users to go beyond
523 ribosomal proteins and leverage other gene families tailored for specific analyses which can
524 improve phylogenetics and the detection of specific taxa. For example, RNAPolA and RNAPolB
525 have been leveraged for phylogeny-guided binning leading to the discovery of missing branches
526 in viral evolution (Gaïa et al. 2023). Furthermore, phylogeography of functional protein families
527 can be leveraged as proxies for microbial metabolism, e.g. phylogeography of ABC transporters
528 can aid in modeling cryptic fluxes of microbial metabolites (Schroer 2023). The EcoPhylo workflow
529 provides a platform for future microbiome projects to benchmark their genome recovery rates
530 upon release of genome collections. Ribosomal protein phylogeography in tandem with reporting
531 read recruitment percentages to representative genome collections, provides comprehensive
532 insights into genome recovery rates given the biodiversity detected in metagenomes. Future
533 studies can leverage the strategy implemented in EcoPhylo to reanalyze existing metagenomic
534 assemblies to identify missing clades or develop tailored methods to optimize overall genome
535 recovery efforts by taking advantage of the increasing availability of genomes and metagenomes.

536 Materials and Methods

537 The EcoPhylo workflow

538 EcoPhylo is a computational workflow implemented in the open-source software ecosystem `anvi'o`
539 (Eren et al. 2015, 2021) using the Python programming language and the workflow management
540 system, Snakemake (Köster and Rahmann 2012). The primary purpose of EcoPhylo is to offer
541 an integrated means to study phylogenetic relationships and ecological distribution patterns of
542 sequences that match to any gene family based on user-provided hidden Markov model (HMM)
543 searches from genomic and metagenomic assemblies. A minimal command line instruction to
544 start an EcoPhylo run is ``anvi-run-workflow -w ecophylo -c config.json``, where ``anvi-run-workflow``
545 is a program in `anvi'o` that runs various workflows, and ``config.json`` is a JSON formatted
546 configuration file that describes file paths (such as the locations of genomes and/or
547 metagenomes) and other parameters (such as the HMM to be used for a homology search, and
548 sequence identity cutoffs). Comprehensive user documentation for EcoPhylo is available at
549 <https://anvio.org/m/ecophylo>.

550 The minimum input for the EcoPhylo is a gene family hidden Markov model (HMM) and a dataset
551 of genomic and/or metagenomic assemblies. EcoPhylo identifies and clusters target genes or
552 translated proteins across assemblies to yield a non-redundant, representative set of open
553 reading frames (ORFs). Next, an amino acid phylogenetic tree is calculated with the translated
554 representative ORFs yielding the evolutionary history captured by homologues from input
555 assemblies. An additional user input to the workflow is a metagenomic sequencing dataset
556 representing ecological sampling or an experimental setup. With this input, the workflow performs
557 metagenomic read recruitment against the representative ORFs to yield ecological insights into

558 the gene family. Finally, the separate data types are integrated into a phylogeographic
559 representation of the gene family (Figure 1).

560 The resulting sequences from the workflow can be organized in the EcoPhylo interactive interface
561 either using an amino acid phylogenetic tree or using hierarchical clustering based on differential
562 read recruitment coverage across metagenomic samples. Additionally, metagenomes can be
563 hierarchically clustered based on the detection of the target gene family. It is recommended to
564 employ hierarchical clustering of metagenomes or sequences in the EcoPhylo interactive
565 interface with the detection read recruitment statistic (rather than coverage values) to minimize
566 the effect of non-specific read recruitment (<https://merenlab.org/anvio-views/>).

567 An application of the EcoPhylo workflow with default settings will (1) identify gene families with
568 the program `hmmsearch` in (Eddy 2011) using the user-provided HMM model, (2) annotate
569 affiliate hmm-hits with taxonomic names with `anvi-run-scg-taxonomy` when applicable, (3)
570 remove hmm-hits with less than 80% HMM model alignment coverage and incomplete ORFs with
571 the anvi'o program `anvi-script-filter-hmm-hits` with parameters `--min-model-coverage 0.8` and
572 `--filter-out-partial-gene-calls` to minimize the inclusion of non-target sequences and spurious
573 HMM hits, (4) dereplicate the resulting DNA sequences at 97% gANI and pick cluster
574 representatives using MMseqs2 (Steinegger and Söding 2017), (5) use the translated
575 representative sequences to calculate a multiple sequence alignment (MSA) using (Edgar 2004)
576 with the `-maxiters 2` flag (Edgar 2004), trim the alignment by removing columns of the alignment
577 with trimAL with the '-gappyout' flag (Capella-Gutiérrez, Silla-Martínez, and Gabaldón 2009), (6)
578 remove sequences that have more than 50% gaps using the anvi'o program `anvi-script-reformat-
579 fasta`, (7) calculate a phylogenetic tree using FastTree (Price, Dehal, and Arkin 2010) with the
580 flag `-fastest`, (8) perform metagenomic read recruitment analysis and profiling of non-translated
581 representative sequences using the anvi'o metagenomic workflow (Shaiber et al. 2020), which by
582 default relies upon Bowtie2 (Langmead and Salzberg 2012), (9) generate miscellaneous data to

583 annotate the representative sequences including taxonomy with ``anvi-estimate-scg-taxonomy``
584 for ribosomal proteins, cluster size, and sequence length, and finally (10) generate `anvi'o` artifacts
585 that give integrated access to the phylogenetic tree of all representative sequences and read
586 recruitment results that can be visualized using the `anvi'o` interactive interface and/or further
587 processed for specific downstream analyses using any popular data analysis environment such
588 as R and/or Python.

589 The workflow that resulted in the recovery and characterization of ribosomal proteins in our
590 manuscript used the following additional steps: (1) we removed input reference genomes that
591 were not detected in at least one of the input metagenomes above a detection value of 0.9 with
592 their ribosomal protein (we kept all MAGs originating from the samples themselves) to only
593 visualize detected populations, (2) we manually curated the ribosomal protein tree when
594 necessary to remove sequences that appeared to be chimeric and those that formed spurious
595 long branches likely originating from metagenomic assembly artifacts, and/or mitochondrial or
596 plastid genomes (Supplementary information) and recalculated new amino acid phylogenetic
597 trees with curated sequences with ``FastTree`` or `IQTREE` with the parameters ``-m WAG -B 1000``
598 (Minh et al. 2020) and imported the new trees using the program ``anvi-import-items-order``, and
599 finally, (3) we generated additional metadata using in-house Python or R scripts and imported
600 additional metadata using the program ``anvi-import-misc-data`` to decorate trees or
601 metagenomes.

602 Benchmarking EcoPhylo workflow with ribosomal proteins using 603 CAMI synthetic metagenomes

604 We validated the EcoPhylo workflow by benchmarking it against the CAMI synthetic
605 metagenomes (Meyer et al. 2022) to identify nucleotide clustering thresholds of ribosomal gene

606 families to limit non-specific read recruitment while maximizing taxonomic resolution. We applied
607 the EcoPhylo workflow across the three CAMI biome synthetic genomic/metagenomic datasets
608 (Marine, Plant-associated, and Strain-madness). As an initial step, we identified the top five most
609 frequent ribosomal gene families that were detected in single-copy in the associated genomic
610 collections for each synthetic metagenomic dataset. We then conducted a parameter grid search,
611 spanning 95%-100% nucleotide similarity parameter grid search (Meyer et al. 2022). Next, we
612 measured the amount of non-specific read recruitment in each EcoPhylo iteration, i.e. reads with
613 equal mapping scores between their primary and secondary alignments (multi-mapped reads),
614 with the following Samtools command: ``samtools view $sample | grep XS:i | cut -f12-13 | sed`
615 `'s/...:i//g' | awk '$1==$2' | wc -l``. The percentage of non-specific read recruitment was calculated
616 by dividing the number of multi-mapped reads by the total number of reads mapped to the
617 representative dataset. With this, we identified that nucleotide clustering thresholds greater than
618 97% began to show signs of non-specific read recruitment (Supplementary information).

619 After identifying 97% nucleotide identity as the optimal threshold, we measured EcoPhylo's ability
620 to contextualize a genomic collection within metagenomic assemblies by quantifying the amount
621 of genomic ribosomal genes clustering with their associated metagenomic assembly ribosomal
622 gene (Supplementary information). Finally, we benchmarked the Shannon diversity and richness
623 captured by different SCGs within the metagenomes and compared it to other taxonomic profiling
624 tools submitted to CAMI (Meyer et al. 2022). To calculate Shannon diversity and richness values
625 for SCGs processed by EcoPhylo we used the R package `vegan` (Dixon 2003) and `Phyloseq`
626 (McMurdie and Holmes 2013). To calculate the richness and alpha diversity values for CAMI
627 ground truth and other profiling tools we extracted relative abundance for each genera included
628 in the associated biome files made available from CAMI. Shannon diversity for SCGs in the
629 EcoPhylo were calculated with the `anvi'o` coverage statistic: Q2Q3 coverage. Datasets were
630 cleaned and visualized with R packages in `Tidyverse` (Wickham et al. 2019).

631 Genome collections

632 All MAG and SAG datasets were filtered for genomes with 50% completion and 10% redundancy
633 using the single-copy core gene collections in *anvi'o* to meet medium-quality draft status in
634 accordance with the community standards (Bowers et al. 2017). For the human oral cavity
635 analysis, 8,615 human oral isolate genomes were downloaded from HOMD v10.1
636 (<https://www.homd.org/ftp/genomes/NCBI/V10.1/>) (Escapa et al. 2018) and 790 MAGs were
637 downloaded from Shaiber et al. (2020) via (doi:[10.6084/m9.figshare.12217805](https://doi.org/10.6084/m9.figshare.12217805),
638 doi:[10.6084/m9.figshare.12217961](https://doi.org/10.6084/m9.figshare.12217961)).

639 For the Hadza tribe human gut microbiome analysis we followed the data download guidelines
640 shared by Carter et al. (2023) to obtain genomes from doi:[10.5281/zenodo.7782708](https://doi.org/10.5281/zenodo.7782708). Carter et al.
641 (2023) formed clusters at 95% gANI by including additional genomes outside of the MAGs they
642 have reconstructed from the Hadza gut metagenomes. To exclusively analyze microbial genomes
643 affiliated with the Hadza metagenomes, we filtered for cluster representatives with cluster
644 members that contained at least one Hadzda adult or infant MAG which produced 2,437
645 representative Bacterial and Archaeal MAGs.

646 Finally, the surface ocean genomic collection was based on Paoli et al. (2022) and augmented
647 with SAGs (Martínez-Pérez et al. 2022) and isolate genomes for SAR11 and *Prochlorococcus*
648 (Delmont and Eren 2018; Delmont et al. 2019a). When metadata was available, we only used
649 genomes sampled from $x < 30$ meters depth to match the surface ocean metagenomic dataset,
650 otherwise, we retained the genomes. The (Paoli et al. 2022) MAG collection included manually
651 curated MAGs from co-assemblies, which included samples from depths deeper than 30 meters
652 in the deep chlorophyll maximum (Delmont et al. 2019a). The final input surface ocean genome
653 dataset contained 1,474 SAGs, 1,723 isolate genomes, and 7,282 MAGs.

654 Metagenome and metagenomic assembly datasets

655 To explore the phylogeography of ribosomal proteins, we used used 71 tooth and plaque
656 metagenomes from the human oral cavity which were downloaded from the NCBI BioProject
657 PRJNA625082 (Shaiber et al. 2020) along with associated co-assemblies
658 (doi:[10.6084/m9.figshare.12217799](https://doi.org/10.6084/m9.figshare.12217799)). Next, to explore deep sequencing in the human gut
659 microbiome we used 388 metagenomes and assemblies from infant and adult members of the
660 Hadza tribe (doi:[10.5281/zenodo.7782708](https://doi.org/10.5281/zenodo.7782708)) using the FTP links shared in from the file
661 `Supplemental_Table_S1.csv` and NCBI BioProject PRJEB49206 (Carter et al. 2023). Finally, to
662 explore the global surface ocean microbiome, we used 237 surface ocean metagenomes and
663 associated assemblies (<30 meters depth) from NCBI BioProjects PRJEB45951 and
664 PRJEB5245228 (Paoli et al. 2022; Sánchez et al. 2023). All metagenomes and associated
665 assembly accessions can be found at Supplementary Table 1.

666 Preprocessing of genomic and metagenomic assemblies and 667 metagenomic short reads

668 Metagenomic and genomic assemblies were preprocessed with the anvi'o contigs workflow with
669 the program `anvi-run-workflow -w contigs` to predict open-reading frames with Prodigal (V2.6.3)
670 and identify SCGs for taxonomic inference with `anvi-run-scg-taxonomy` (Hyatt et al. 2010;
671 Shaiber et al. 2020). No contig size filters were implemented during this process to include
672 ribosomal proteins located on small contigs. To limit detection of misassemblies in downstream
673 analyses, only ribosomal proteins with complete open-reading frames (as predicted by Prodigal)
674 were analyzed with EcoPhylo (Hyatt et al. 2010). Additionally, metagenomic samples were quality
675 controlled with the anvi'o metagenomics workflow with the program `anvi-run-workflow -w
676 metagenomics` (Shaiber et al. 2020). This workflow uses the tool `iu-filter-quality-minoche` (Eren

677 et al. 2013), which implements methods described in (Minoche, Dohm, and Himmelbauer 2011).
678 All Snakemake workflows in this manuscript leveraged Snakemake v7.32.4 (Köster and Rahmann
679 2012).

680 Gene-level taxonomy of ribosomal proteins

681 To assign gene level taxonomy to ribosomal proteins, the EcoPhylo workflow relies upon the
682 anvi'o tools ``anvi-run-scg-taxonomy`` and ``anvi-estimate-scg-taxonomy``, which leverage the
683 genomes and their taxonomy made available by the GTDB (Parks et al. 2022) to identify
684 taxonomic affiliations of genes that match to any of the ribosomal proteins L1, L13, L14, L16, L17,
685 L19, L2, L20, L21p, L22, L27A, L3, L4, L5, S11, S15, S16, S2, S6, S7, S8, or S9. During the
686 workflow, EcoPhylo uses ``anvi-run-scg-taxonomy`` to search for ribosomal genes annotated within
687 each anvi'o contigs database against the downloaded marker gene dataset with DIAMOND
688 v0.9.14 (Buchfink, Reuter, and Drost 2021). Later in the workflow, EcoPhylo runs ``anvi-estimate-
689 scg-taxonomy --metagenome-mode`` on the representative set of ribosomal proteins, which
690 assigns a consensus taxonomy to each sequence. The program ``anvi-estimate-scg-taxonomy``
691 does not provide a taxonomic annotation if the ribosomal protein is less than 90% similar to any
692 of the ribosomal proteins found in GTDB genomes. In some cases, ribosomal proteins without
693 taxonomic annotation can be manually annotated with taxonomy based on the annotated
694 sequences that surround them in the phylogenetic tree, as we described in the section
695 "Taxonomic binning to improve genome recovery estimations".

696 Selection of ribosomal proteins to contextualize genomic 697 collections in metagenomes

698 To pick ribosomal gene families to study genome collections, we selected ribosomal genes that
699 were annotated in the majority of genomes in single-copy. We then cross-referenced selected
700 ribosomal genes with their assembly rates in metagenomes and disregarded candidate ribosomal
701 gene families that were under- or over-assembled in the dataset. To do this, we ran the EcoPhylo
702 workflow with the input dataset of genomic and metagenomic assemblies until the rule
703 ``process_hmm_hits``, which will filter for high-quality HMM-hits as described above. Finally, we
704 extracted ribosomal protein hits from all assemblies with the `anvi'o` command ``anvi-script-gen-
705 hmm-hits-matrix-across-genomes`` and tabulated/visualized the distribution in R using the
706 Tidyverse (Wickham et al. 2019).

707 Distribution of HMM alignment coverage and SCG detection across 708 GTDB

709 To identify optimal ribosomal proteins and HMM hit filtering thresholds, we explored the
710 distribution of SCG detection and HMM alignment coverage across GTDB genomes. The analysis
711 used the first two rules of the EcoPhylo workflow (`anvi_run_hmms_hmmsearch` and
712 `filter_hmm_hits_by_model_coverage`) to annotate the RefSeq representative genomes from
713 GTDB release 95 (Parks et al. 2020), with the single-copy core gene HMM collections included in
714 `anvi'o`. The first rule of the workflow used the program ``hmmsearch`` in (Eddy 2011) to identify HMM
715 hits, while the second rule was modified to include all HMM hit model coverage values by setting
716 the parameter ``anvi-script-filter-hmm-hits-table --min-model-coverage 0``. We stopped the
717 workflow after this rule and visualized the raw distribution of model and gene coverage values

718 from `hmmsearch --domtblout` output file leading us to to identify an 80% HMM hit model
719 coverage as an optimal filtering threshold to identify ribosomal proteins. Next, we restarted the
720 workflow but re-modified the second rule parameter `anvi-script-filter-hmm-hits-table --min-model-
721 coverage 0.8` to filter for HMMs hits with at least 80% model alignment coverage. Finally, we
722 extracted all ribosomal gene families from the genome dataset with anvi'o program `anvi-script-
723 gen-hmm-hits-matrix-across-genomes` and visualized the genome detection and SCG copy
724 number across the dataset in in R using the Tidyverse (Wickham et al. 2019).

725 Detection of whole genomes in metagenomic data

726 In some cases, ribosomal proteins clustering at 97% brought together large groups of highly
727 similar isolate genomes. To identify the specific genome that is detected in the metagenomic
728 datasets, we re-clustered the target EcoPhylo protein at 98% to resolve sequence clusters and
729 thus increase the number of representative sequences. We then used the whole genomes
730 associated with the new, larger set of representative proteins to explore their distribution in
731 metagenomes by performing the anvi'o metagenomic workflow (Shaiber et al. 2020). Our
732 threshold for detection of a whole-genome in metagenomic data was 50% (percent of genome
733 covered by at least one read from metagenomic read recruitment), which was found to be efficient
734 for human oral cavity microbes (Utter et al. 2020).

735 Genome recovery rate estimations

736 Genome recovery rates were estimated to measure which individual or combination of genome
737 types (MAGs, SAGs, isolate genomes) most effectively sampled clades in the ribosomal protein
738 phylogenetic trees calculated during the EcoPhylo workflow. To calculate genome recovery rates
739 for any given taxon, we divided the number of sequence clusters that contained a sequence from

740 a given genome recovery method to the total number of representative sequences EcoPhylo
741 reported for that taxon. Taxonomic assignments of sequence cluster representatives were
742 determined with `anvi-estimate-scg-taxonomy`.

743 Taxonomic binning to improve genome recovery estimations

744 A subset of ribosomal proteins lacked taxon assignments from `anvi-estimate-scg-taxonomy` due
745 to their sequence similarity being $x < 90\%$ to GTDB genomes (See methods section: Gene-level
746 taxonomy of ribosomal proteins). Using the `anvi-interactive` interface, we examined the
747 placement of these proteins in the EcoPhylo ribosomal protein phylogenetic tree and manually
748 assigned taxon names based on the taxonomic affiliations of neighboring sequences.
749 Unannotated sequences were assigned taxonomy only when phylogenetic clustering
750 demonstrated clear consistency among neighboring sequences. These refined taxonomic
751 annotations were used to improve estimations of genome recovery in the main figures (*rpL19* and
752 *rpS15* in Figure 2 and *rpL14* in Figure 3).

753 Data and code availability

754 The URL <https://merenlab.org/data/ecophylo-ribosomal-proteins/> serves all code and data
755 needed to reproduce our study. Additionally, all anvi'o artifacts that give interactive access to
756 EcoPhylo interfaces are publicly available at doi:[10.6084/m9.figshare.28207481](https://doi.org/10.6084/m9.figshare.28207481). Publicly
757 available genomes and metagenomes we used in our study are listed in the Supplementary
758 Tables, which are available via doi:[10.6084/m9.figshare.28200050](https://doi.org/10.6084/m9.figshare.28200050), along with the Supplementary
759 Information text.

760 Acknowledgements

761 We thank all authors who made the raw metagenomic reads, metagenomic assemblies, and
762 genomes from their manuscripts publically available and conveniently accessible for secondary
763 analyses. We also thank the members of the Meren Lab (<https://merenlab.org/people/>), the Light
764 Lab (<https://www.lightlab.uchicago.edu/people/>), and the Blekman Lab
765 (<http://blekmanlab.org/members.html>) for helpful discussions and whiteboard sessions. We are
766 also thankful to Pedram Esfahani, Kimberly Gräsch, and the rest of the University of Chicago
767 Center for Research and Computing Center for their patience and support. MSS acknowledges
768 support from NIH Genetics and Regulation Training Grant (T32 GM07197). AME acknowledges
769 support from the Center for Chemical Currencies of a Microbial Planet (C-CoMP) (NSF Award
770 OCE-2019589, C-CoMP publication #070), and Simons Foundation (grant #687269).

771 Author contributions

772 MSS and AME conceptualized the study. MSS curated data and performed formal analyses.
773 MSS, IAV, MLK, MS, SEM, and AME developed software tools. MSS, FT, and AME interpreted
774 findings. LM, TOD, and SHL helped with interpretation of results. MSS and AME wrote the original
775 draft of the study. SHL and AME managed the project and acquired funding. All authors
776 commented on and made suggestions, and approved the final manuscript.

777 Competing interests

778 Authors declare that they have no conflicts of interest.

779 References

- 780 Al-Shayeb, Basem, Petr Skopintsev, Katarzyna M. Soczek, Elizabeth C. Stahl, Zheng Li, Evan
781 Groover, Dylan Smock, et al. 2022. "Diverse Virus-Encoded CRISPR-Cas Systems Include
782 Streamlined Genome Editors." *Cell* 185 (24): 4574–86.e16.
- 783 Baker-Austin, Craig, Joaquin Trinanes, Narjol Gonzalez-Escalona, and Jaime Martinez-Urtaza.
784 2017. "Non-Cholera Vibrios: The Microbial Barometer of Climate Change." *Trends in*
785 *Microbiology* 25 (1): 76–84.
- 786 Biller, Steven J., Paul M. Berube, Keven Dooley, Madeline Williams, Brandon M. Satinsky,
787 Thomas Hackl, Shane L. Hogle, et al. 2018. "Marine Microbial Metagenomes Sampled
788 across Space and Time." *Scientific Data* 5 (September):180176.
- 789 Biller, Steven J., Paul M. Berube, Debbie Lindell, and Sallie W. Chisholm. 2015.
790 "Prochlorococcus: The Structure and Function of Collective Diversity." *Nature Reviews.*
791 *Microbiology* 13 (1): 13–27.
- 792 Bowen, William H., Robert A. Burne, Hui Wu, and Hyun Koo. 2018. "Oral Biofilms: Pathogens,
793 Matrix, and Polymicrobial Interactions in Microenvironments." *Trends in Microbiology* 26 (3):
794 229–42.
- 795 Bowers, Robert M., Nikos C. Kyrpides, Ramunas Stepanauskas, Miranda Harmon-Smith, Devin
796 Doud, T. B. K. Reddy, Frederik Schulz, et al. 2017. "Minimum Information about a Single
797 Amplified Genome (MISAG) and a Metagenome-Assembled Genome (MIMAG) of Bacteria
798 and Archaea." *Nature Biotechnology* 35 (8): 725–31.
- 799 Brown, Christopher T., Laura A. Hug, Brian C. Thomas, Itai Sharon, Cindy J. Castelle, Andrea
800 Singh, Michael J. Wilkins, Kelly C. Wrighton, Kenneth H. Williams, and Jillian F. Banfield.
801 2015. "Unusual Biology across a Group Comprising More than 15% of Domain Bacteria."
802 *Nature* 523 (7559): 208–11.
- 803 Brown, Mark V., Federico M. Lauro, Matthew Z. DeMaere, Les Muir, David Wilkins, Torsten
804 Thomas, Martin J. Riddle, et al. 2012. "Global Biogeography of SAR11 Marine Bacteria."
805 *Molecular Systems Biology* 8 (1): 595.
- 806 Buchfink, Benjamin, Klaus Reuter, and Hajk-Georg Drost. 2021. "Sensitive Protein Alignments
807 at Tree-of-Life Scale Using DIAMOND." *Nature Methods* 18 (4): 366–68.
- 808 Capella-Gutiérrez, Salvador, José M. Silla-Martínez, and Toni Gabaldón. 2009. "trimAl: A Tool
809 for Automated Alignment Trimming in Large-Scale Phylogenetic Analyses." *Bioinformatics*
810 25 (15): 1972–73.
- 811 Carter, Matthew M., Matthew R. Olm, Bryan D. Merrill, Dylan Dahan, Surya Tripathi, Sean P.
812 Spencer, Feiqiao B. Yu, et al. 2023. "Ultra-Deep Sequencing of Hadza Hunter-Gatherers
813 Recovers Vanishing Gut Microbes." *Cell*, June. <https://doi.org/10.1016/j.cell.2023.05.046>.
- 814 Chang, Tianyi, Gregory S. Gavelis, Julia M. Brown, and Ramunas Stepanauskas. 2024.
815 "Genomic Representativeness and Chimerism in Large Collections of SAGs and MAGs of
816 Marine Prokaryoplankton." *Microbiome* 12 (1): 126.
- 817 Chen, Jianwei, Yangyang Jia, Ying Sun, Kun Liu, Changhao Zhou, Chuan Liu, Denghui Li, et al.
818 2024a. "Global Marine Microbial Diversity and Its Potential in Bioprospecting." *Nature* 633
819 (8029): 371–79.
- 820 ———. 2024b. "Global Marine Microbial Diversity and Its Potential in Bioprospecting." *Nature*,
821 September. <https://doi.org/10.1038/s41586-024-07891-2>.
- 822 Chen, Lin-Xing, Karthik Anantharaman, Alon Shaiber, A. Murat Eren, and Jillian F. Banfield.
823 2020. "Accurate and Complete Genomes from Metagenomes." *Genome Research* 30 (3):
824 315–33.
- 825 Chen, Tsute, Wen-Han Yu, Jacques Izard, Oxana V. Baranova, Abirami Lakshmanan, and

- 826 Floyd E. Dewhirst. 2010. "The Human Oral Microbiome Database: A Web Accessible
827 Resource for Investigating Oral Microbe Taxonomic and Genomic Information." *Database:
828 The Journal of Biological Databases and Curation* 2010 (0): baq013.
- 829 Comman, Andre, Jacob West-Roberts, Antonio Pedro Camargo, Simon Roux, Martin
830 Beracochea, Milot Mirdita, Sergey Ovchinnikov, and Yunha Hwang. 2024. "The OMG
831 Dataset: An Open MetaGenomic Corpus for Mixed-Modality Genomic Language Modeling."
832 *bioRxiv*. <https://doi.org/10.1101/2024.08.14.607850>.
- 833 Crits-Christoph, Alexander, Spencer Diamond, Basem Al-Shayeb, Luis Valentin-Alvarado, and
834 Jillian F. Banfield. 2022. "A Widely Distributed Genus of Soil Acidobacteria Genomically
835 Enriched in Biosynthetic Gene Clusters." *ISME Communications* 2 (1): 1–8.
- 836 Cross, Karissa L., James H. Campbell, Manasi Balachandran, Alisha G. Campbell, Connor J.
837 Cooper, Ann Griffen, Matthew Heaton, et al. 2019. "Targeted Isolation and Cultivation of
838 Uncultivated Bacteria by Reverse Genomics." *Nature Biotechnology* 37 (11): 1314–21.
- 839 Delmont, Tom O., and A. Murat Eren. 2018. "Linking Pangenomes and Metagenomes: The
840 Prochlorococcus Metapangenome." *PeerJ* 6 (January): e4320.
- 841 Delmont, Tom O., Evan Kiefl, Ozsel Kilinc, Ozcan C. Esen, Ismail Uysal, Michael S. Rappé,
842 Steven Giovannoni, and A. Murat Eren. 2019a. "Single-Amino Acid Variants Reveal
843 Evolutionary Processes That Shape the Biogeography of a Global SAR11 Subclade." *eLife*
844 8 (September). <https://doi.org/10.7554/eLife.46497>.
- 845 ———. 2019b. "Single-Amino Acid Variants Reveal Evolutionary Processes That Shape the
846 Biogeography of a Global SAR11 Subclade." *eLife* 8 (September).
847 <https://doi.org/10.7554/eLife.46497>.
- 848 Delmont, Tom O., Christopher Quince, Alon Shaiber, Özcan C. Esen, Sonny Tm Lee, Michael
849 S. Rappé, Sandra L. McLellan, Sebastian Lücker, and A. Murat Eren. 2018a. "Nitrogen-
850 Fixing Populations of Planctomycetes and Proteobacteria Are Abundant in Surface Ocean
851 Metagenomes." *Nature Microbiology* 3 (7): 804–13.
- 852 ———. 2018b. "Nitrogen-Fixing Populations of Planctomycetes and Proteobacteria Are
853 Abundant in Surface Ocean Metagenomes." *Nature Microbiology* 3 (7): 804–13.
- 854 Dewhirst, Floyd E., Tuste Chen, Jacques Izard, Bruce J. Paster, Anne C. R. Tanner, Wen-Han
855 Yu, Abirami Lakshmanan, and William G. Wade. 2010. "The Human Oral Microbiome."
856 *Journal of Bacteriology* 192 (19): 5002–17.
- 857 Diamond, Spencer, Peter F. Andeer, Zhou Li, Alexander Crits-Christoph, David Burstein, Karthik
858 Anantharaman, Katherine R. Lane, et al. 2019. "Mediterranean Grassland Soil C-N
859 Compound Turnover Is Dependent on Rainfall and Depth, and Is Mediated by Genomically
860 Divergent Microorganisms." *Nature Microbiology* 4 (8): 1356–67.
- 861 Dixon, Philip. 2003. "VEGAN, a Package of R Functions for Community Ecology." *Journal of
862 Vegetation Science: Official Organ of the International Association for Vegetation Science*
863 14 (6): 927–30.
- 864 Durrant, Matthew G., Alison Fanton, Josh Tycko, Michaela Hinks, Sita S. Chandrasekaran,
865 Nicholas T. Perry, Julia Schaepe, et al. 2023. "Systematic Discovery of Recombinases for
866 Efficient Integration of Large DNA Sequences into the Human Genome." *Nature
867 Biotechnology* 41 (4): 488–99.
- 868 Eddy, Sean R. 2011. "Accelerated Profile HMM Searches." *PLoS Computational Biology* 7 (10):
869 e1002195.
- 870 Edgar, Robert C. 2004. "MUSCLE: Multiple Sequence Alignment with High Accuracy and High
871 Throughput." *Nucleic Acids Research* 32 (5): 1792–97.
- 872 Emerson, Joanne B., Brian C. Thomas, Walter Alvarez, and Jillian F. Banfield. 2016.
873 "Metagenomic Analysis of a High Carbon Dioxide Subsurface Microbial Community
874 Populated by Chemolithoautotrophs and Bacteria and Archaea from Candidate Phyla."
875 *Environmental Microbiology* 18 (6): 1686–1703.
- 876 Eren, A. Murat, and Jillian F. Banfield. 2024. "Modern Microbiology: Embracing Complexity

- 877 through Integration across Scales.” *Cell* 187 (19): 5151–70.
- 878 Eren, A. Murat, Gary G. Borisy, Susan M. Huse, and Jessica L. Mark Welch. 2014. “Oligotyping
879 Analysis of the Human Oral Microbiome.” *Proceedings of the National Academy of
880 Sciences of the United States of America* 111 (28): E2875–84.
- 881 Eren, A. Murat, Özcan C. Esen, Christopher Quince, Joseph H. Vineis, Hilary G. Morrison,
882 Mitchell L. Sogin, and Tom O. Delmont. 2015. “Anvi’o: An Advanced Analysis and
883 Visualization Platform for ‘Omics Data.” *PeerJ* 3 (October):e1319.
- 884 Eren, A. Murat, Evan Kiefl, Alon Shaiber, Iva Veseli, Samuel E. Miller, Matthew S. Schechter,
885 Isaac Fink, et al. 2021. “Community-Led, Integrated, Reproducible Multi-Omics with Anvi’o.”
886 *Nature Microbiology* 6 (1): 3–6.
- 887 Eren, A. Murat, Joseph H. Vineis, Hilary G. Morrison, and Mitchell L. Sogin. 2013. “A Filtering
888 Method to Generate High Quality Short Reads Using Illumina Paired-End Technology.”
889 *PLoS One* 8 (6): e66643.
- 890 Escapa, Isabel F., Tsute Chen, Yanmei Huang, Prasad Gajare, Floyd E. Dewhirst, and
891 Katherine P. Lemon. 2018. “New Insights into Human Nostril Microbiome from the
892 Expanded Human Oral Microbiome Database (eHOMD): A Resource for the Microbiome of
893 the Human Aerodigestive Tract.” *mSystems* 3 (6). [https://doi.org/10.1128/mSystems.00187-](https://doi.org/10.1128/mSystems.00187-18)
894 18.
- 895 Freel, Kelle C., Sarah J. Tucker, Evan B. Freel, Stephen J. Giovannoni, A. Murat Eren, and
896 Michael S. Rappe. 2024. “New Isolate Genomes and Global Marine Metagenomes Resolve
897 Ecologically Relevant Units of SAR11.” *Microbiology*. bioRxiv.
898 <https://www.biorxiv.org/content/10.1101/2024.12.24.630191v2>.
- 899 Gaïa, Morgan, Lingjie Meng, Eric Pelletier, Patrick Forterre, Chiara Vanni, Antonio Fernandez-
900 Guerra, Olivier Jaillon, et al. 2023. “Mirusviruses Link Herpesviruses to Giant Viruses.”
901 *Nature* 616 (7958): 783–89.
- 902 Giovannoni, Stephen J. 2017. “SAR11 Bacteria: The Most Abundant Plankton in the Oceans.”
903 *Annual Review of Marine Science* 9 (January):231–55.
- 904 Halloran, Kathryn H., Rogier Braakman, Allison Coe, Gretchen Swarr, Melissa C. Kido Soule,
905 Sallie W. Chisholm, and Elizabeth B. Kujawinski. 2025. “Uptake of Prochlorococcus-
906 Derived Metabolites by *Alteromonas Macleodii* MIT1002 Shows High Levels of Substrate
907 Specificity.” *Microbiology*. bioRxiv.
908 <https://www.biorxiv.org/content/10.1101/2025.01.10.632383v1>.
- 909 Hamilton, Trinity L., Roderick J. Bovee, Sarah R. Sattin, Wiebke Mohr, William P. Gilhooly 3rd,
910 Timothy W. Lyons, Ann Pearson, and Jennifer L. Macalady. 2016. “Carbon and Sulfur
911 Cycling below the Chemocline in a Meromictic Lake and the Identification of a Novel
912 Taxonomic Lineage in the FCB Superphylum, Candidatus Aegiribacteria.” *Frontiers in
913 Microbiology* 7 (April):598.
- 914 Henríquez-Castillo, Carlos, Alvaro M. Plominsky, Salvador Ramírez-Flandes, Anthony D.
915 Bertagnolli, Frank J. Stewart, and Osvaldo Ulloa. 2022. “Metaomics Unveils the
916 Contribution of *Alteromonas* Bacteria to Carbon Cycling in Marine Oxygen Minimum
917 Zones.” *Frontiers in Marine Science* 9 (September).
918 <https://doi.org/10.3389/fmars.2022.993667>.
- 919 He, Xuesong, Jeffrey S. McLean, Anna Edlund, Shibu Yooseph, Adam P. Hall, Su-Yang Liu,
920 Pieter C. Dorrestein, et al. 2015. “Cultivation of a Human-Associated TM7 Phylotype
921 Reveals a Reduced Genome and Epibiotic Parasitic Lifestyle.” *Proceedings of the National
922 Academy of Sciences of the United States of America* 112 (1): 244–49.
- 923 Hosokawa, Masahito, and Yohei Nishikawa. 2024. “Tools for Microbial Single-Cell Genomics for
924 Obtaining Uncultured Microbial Genomes.” *Biophysical Reviews* 16 (1): 69–77.
- 925 Huang, Sijun, Steven W. Wilhelm, H. Rodger Harvey, Karen Taylor, Nianzhi Jiao, and Feng
926 Chen. 2012. “Novel Lineages of *Prochlorococcus* and *Synechococcus* in the Global
927 Oceans.” *The ISME Journal* 6 (2): 285–97.

- 928 Hug, Laura A., Cindy J. Castelle, Kelly C. Wrighton, Brian C. Thomas, Itai Sharon, Kyle R.
929 Frischkorn, Kenneth H. Williams, Susannah G. Tringe, and Jillian F. Banfield. 2013.
930 “Community Genomic Analyses Constrain the Distribution of Metabolic Traits across the
931 Chloroflexi Phylum and Indicate Roles in Sediment Carbon Cycling.” *Microbiome* 1 (1): 22.
932 Hwang, Yunha, Andre L. Cornman, Elizabeth H. Kellogg, Sergey Ovchinnikov, and Peter R.
933 Girguis. 2024. “Genomic Language Model Predicts Protein Co-Regulation and Function.”
934 *Nature Communications* 15 (1): 2880.
- 935 Hyatt, Doug, Gwo-Liang Chen, Philip F. Locascio, Miriam L. Land, Frank W. Larimer, and Loren
936 J. Hauser. 2010. “Prodigal: Prokaryotic Gene Recognition and Translation Initiation Site
937 Identification.” *BMC Bioinformatics* 11 (March):119.
- 938 Imachi, Hiroyuki, Masaru K. Nobu, Nozomi Nakahara, Yuki Morono, Miyuki Ogawara, Yoshihiro
939 Takaki, Yoshinori Takano, et al. 2020. “Isolation of an Archaeon at the Prokaryote–
940 eukaryote Interface.” *Nature*. <https://doi.org/10.1038/s41586-019-1916-6>.
- 941 Jain, Chirag, Luis M. Rodriguez-R, Adam M. Phillippy, Konstantinos T. Konstantinidis, and
942 Srinivas Aluru. 2018. “High Throughput ANI Analysis of 90K Prokaryotic Genomes Reveals
943 Clear Species Boundaries.” *Nature Communications* 9 (1): 5114.
- 944 Jiang, Cheng-Ying, Libing Dong, Jian-Kang Zhao, Xiaofang Hu, Chaohua Shen, Yuxin Qiao,
945 Xinyue Zhang, et al. 2016. “High-Throughput Single-Cell Cultivation on Microfluidic Streak
946 Plates.” *Applied and Environmental Microbiology* 82 (7): 2210–18.
- 947 Johnson, Zackary I., Erik R. Zinser, Allison Coe, Nathan P. McNulty, E. Malcolm S. Woodward,
948 and Sallie W. Chisholm. 2006. “Niche Partitioning among Prochlorococcus Ecotypes along
949 Ocean-Scale Environmental Gradients.” *Science (New York, N.Y.)* 311 (5768): 1737–40.
- 950 Kauffman, Kathryn M., Fatima A. Hussain, Joy Yang, Philip Arevalo, Julia M. Brown, William K.
951 Chang, David VanInsberghe, et al. 2018. “A Major Lineage of Non-Tailed dsDNA Viruses
952 as Unrecognized Killers of Marine Bacteria.” *Nature* 554 (7690): 118–22.
- 953 Kent, Alyssa G., Steven E. Baer, Céline Mougnot, Jeremy S. Huang, Alyse A. Larkin, Michael
954 W. Lomas, and Adam C. Martiny. 2019. “Parallel Phylogeography of Prochlorococcus and
955 Synechococcus.” *The ISME Journal* 13 (2): 430–41.
- 956 Kent, Alyssa G., Chris L. Dupont, Shibu Yooseph, and Adam C. Martiny. 2016. “Global
957 Biogeography of Prochlorococcus Genome Diversity in the Surface Ocean.” *The ISME*
958 *Journal* 10 (8): 1856–65.
- 959 Kessel, Julia C. van, and Andrew Camilli. 2024. “Vibrio Cholerae: A Fundamental Model System
960 for Bacterial Genetics and Pathogenesis Research.” *Journal of Bacteriology* 206 (11):
961 e0024824.
- 962 Köster, Johannes, and Sven Rahmann. 2012. “Snakemake--a Scalable Bioinformatics Workflow
963 Engine.” *Bioinformatics (Oxford, England)* 28 (19): 2520–22.
- 964 Langmead, Ben, and Steven L. Salzberg. 2012. “Fast Gapped-Read Alignment with Bowtie 2.”
965 *Nature Methods* 9 (4): 357–59.
- 966 Lu, Jennifer, Florian P. Breitwieser, Peter Thielen, and Steven L. Salzberg. 2017. “Bracken:
967 Estimating Species Abundance in Metagenomics Data.” *PeerJ. Computer Science* 3
968 (e104): e104.
- 969 Lu, Zhiying, Elizabeth Entwistle, Matthew D. Kuhl, Alexander R. Durrant, Marcelo Malisano
970 Barreto Filho, Anuradha Goswami, and J. Jeffrey Morris. 2024. “Coevolution of Marine
971 Phytoplankton and Alteromonas Bacteria in Response to pCO₂ and Co-Culture.” *The ISME*
972 *Journal*, December, wrae259.
- 973 Ma, Bin, Caiyu Lu, Yiling Wang, Jingwen Yu, Kankan Zhao, Ran Xue, Hao Ren, et al. 2023. “A
974 Genomic Catalogue of Soil Microbiomes Boosts Mining of Biodiversity and Genetic
975 Resources.” *Nature Communications* 14 (1): 7318.
- 976 Malmstrom, Rex R., Sébastien Rodrigue, Katherine H. Huang, Libusha Kelly, Suzanne E. Kern,
977 Anne Thompson, Sara Roggensack, Paul M. Berube, Matthew R. Henn, and Sallie W.
978 Chisholm. 2013. “Ecology of Uncultured Prochlorococcus Clades Revealed through Single-

- 979 Cell Genomics and Biogeographic Analysis.” *The ISME Journal* 7 (1): 184–98.
- 980 Manck, Lauren E., Jiwoon Park, Benjamin J. Tully, Alfonso M. Poire, Randelle M. Bundy,
981 Christopher L. Dupont, and Katherine A. Barbeau. 2022. “Petrobactin, a Siderophore
982 Produced by *Alteromonas*, Mediates Community Iron Acquisition in the Global Ocean.” *The*
983 *ISME Journal* 16 (2): 358–69.
- 984 Manghi, Paolo, Aitor Blanco-Míguez, Serena Manara, Amir NabiNejad, Fabio Cumbo,
985 Francesco Beghini, Federica Armanini, et al. 2023. “MetaPhlAn 4 Profiling of Unknown
986 Species-Level Genome Bins Improves the Characterization of Diet-Associated Microbiome
987 Changes in Mice.” *Cell Reports* 42 (5): 112464.
- 988 Mark Welch, Jessica L., Floyd E. Dewhirst, and Gary G. Borisy. 2019. “Biogeography of the Oral
989 Microbiome: The Site-Specialist Hypothesis.” *Annual Review of Microbiology* 73
990 (September):335–58.
- 991 Martínez-Pérez, Clara, Chris Greening, Sean K. Bay, Rachael J. Lappan, Zihao Zhao, Daniele
992 De Corte, Christina Hulbe, et al. 2022. “Phylogenetically and Functionally Diverse
993 Microorganisms Reside under the Ross Ice Shelf.” *Nature Communications* 13 (1): 117.
- 994 Matheus Carnevali, Paula B., Adi Lavy, Alex D. Thomas, Alexander Crits-Christoph, Spencer
995 Diamond, Raphaël Méheust, Matthew R. Olm, et al. 2021. “Meanders as a Scaling Motif for
996 Understanding of Floodplain Soil Microbiome and Biogeochemical Potential at the
997 Watershed Scale.” *Microbiome* 9 (1): 121.
- 998 McMurdie, Paul J., and Susan Holmes. 2013. “Phyloseq: An R Package for Reproducible
999 Interactive Analysis and Graphics of Microbiome Census Data.” *PloS One* 8 (4): e61217.
- 1000 Meyer, Fernando, Adrian Fritz, Zhi-Luo Deng, David Koslicki, Till Robin Lesker, Alexey
1001 Gurevich, Gary Robertson, et al. 2022. “Critical Assessment of Metagenome Interpretation:
1002 The Second Round of Challenges.” *Nature Methods* 19 (4): 429–40.
- 1003 Minh, Bui Quang, Heiko A. Schmidt, Olga Chernomor, Dominik Schrempf, Michael D.
1004 Woodhams, Arndt von Haeseler, and Robert Lanfear. 2020. “IQ-TREE 2: New Models and
1005 Efficient Methods for Phylogenetic Inference in the Genomic Era.” *Molecular Biology and*
1006 *Evolution* 37 (5): 1530–34.
- 1007 Minoche, André E., Juliane C. Dohm, and Heinz Himmelbauer. 2011. “Evaluation of Genomic
1008 High-Throughput Sequencing Data Generated on Illumina HiSeq and Genome Analyzer
1009 Systems.” *Genome Biology* 12 (11): R112.
- 1010 Morris, Robert M., Michael S. Rappé, Stephanie A. Connon, Kevin L. Vergin, William A. Siebold,
1011 Craig A. Carlson, and Stephen J. Giovannoni. 2002. “SAR11 Clade Dominates Ocean
1012 Surface Bacterioplankton Communities.” *Nature* 420 (6917): 806–10.
- 1013 Naether, Daniela J., Slavtsho Slawtschew, Sebastian Stasik, Maria Engel, Martin Olzog, Lukas
1014 Y. Wick, Kenneth N. Timmis, and Hermann J. Heipieper. 2013. “Adaptation of the
1015 Hydrocarbonoclastic Bacterium *Alcanivorax Borkumensis* SK2 to Alkanes and Toxic
1016 Organic Compounds: A Physiological and Transcriptomic Approach.” *Applied and*
1017 *Environmental Microbiology* 79 (14): 4282–93.
- 1018 Nayfach, Stephen, Beltran Rodriguez-Mueller, Nandita Garud, and Katherine S. Pollard. 2016.
1019 “An Integrated Metagenomics Pipeline for Strain Profiling Reveals Novel Patterns of
1020 Bacterial Transmission and Biogeography.” *Genome Research* 26 (11): 1612–25.
- 1021 Nguyen, Eric, Michael Poli, Matthew G. Durrant, Armin W. Thomas, Brian Kang, Jeremy
1022 Sullivan, Madelena Y. Ng, et al. 2024. “Sequence Modeling and Design from Molecular to
1023 Genome Scale with Evo.” <https://doi.org/10.1101/2024.02.27.582234>.
- 1024 Olm, Matthew R., Alexander Crits-Christoph, Spencer Diamond, Adi Lavy, Paula B. Matheus
1025 Carnevali, and Jillian F. Banfield. 2020. “Consistent Metagenome-Derived Metrics Verify
1026 and Delineate Bacterial Species Boundaries.” *mSystems* 5 (1).
1027 <https://doi.org/10.1128/mSystems.00731-19>.
- 1028 Pachiadaki, Maria G., Julia M. Brown, Joseph Brown, Oliver Bezuidt, Paul M. Berube, Steven J.
1029 Biller, Nicole J. Poulton, et al. 2019. “Charting the Complexity of the Marine Microbiome

- 1030 through Single-Cell Genomics.” *Cell* 179 (7): 1623–35.e11.
- 1031 Paoli, Lucas, Hans-Joachim Ruscheweyh, Clarissa C. Forneris, Florian Hubrich, Satria Kautsar,
1032 Agneya Bhushan, Alessandro Lotti, et al. 2022. “Biosynthetic Potential of the Global Ocean
1033 Microbiome.” *Nature* 607 (7917): 111–18.
- 1034 Parks, Donovan H., Maria Chuvochina, Pierre-Alain Chaumeil, Christian Rinke, Aaron J.
1035 Mussig, and Philip Hugenholtz. 2020. “A Complete Domain-to-Species Taxonomy for
1036 Bacteria and Archaea.” *Nature Biotechnology* 38 (9): 1079–86.
- 1037 Parks, Donovan H., Maria Chuvochina, Christian Rinke, Aaron J. Mussig, Pierre-Alain
1038 Chaumeil, and Philip Hugenholtz. 2022. “GTDB: An Ongoing Census of Bacterial and
1039 Archaeal Diversity through a Phylogenetically Consistent, Rank Normalized and Complete
1040 Genome-Based Taxonomy.” *Nucleic Acids Research* 50 (D1): D785–94.
- 1041 Parks, Donovan H., Christian Rinke, Maria Chuvochina, Pierre-Alain Chaumeil, Ben J.
1042 Woodcroft, Paul N. Evans, Philip Hugenholtz, and Gene W. Tyson. 2017. “Recovery of
1043 Nearly 8,000 Metagenome-Assembled Genomes Substantially Expands the Tree of Life.”
1044 *Nature Microbiology* 2 (11): 1533–42.
- 1045 Pasolli, Edoardo, Francesco Asnicar, Serena Manara, Moreno Zolfo, Nicolai Karcher, Federica
1046 Armanini, Francesco Beghini, et al. 2019. “Extensive Unexplored Human Microbiome
1047 Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age,
1048 Geography, and Lifestyle.” *Cell* 176 (3): 649–62.e20.
- 1049 Pedler, Byron E., Lihini I. Aluwihare, and Farooq Azam. 2014. “Single Bacterial Strain Capable
1050 of Significant Contribution to Carbon Cycling in the Surface Ocean.” *Proceedings of the
1051 National Academy of Sciences of the United States of America* 111 (20): 7202–7.
- 1052 Prasad, Manoj, Nozomu Obana, Kaori Sakai, Toshiki Nagakubo, Shun Miyazaki, Masanori
1053 Toyofuku, Jacques Fattaccioli, Nobuhiko Nomura, and Andrew S. Utada. 2019. “Point
1054 Mutations Lead to Increased Levels of c-Di-GMP and Phenotypic Changes to the Colony
1055 Biofilm Morphology in *Alcanivorax Borkumensis* SK2.” *Microbes and Environments* 34 (1):
1056 104–7.
- 1057 Prasad, M., N. Obana, S-Z Lin, S. Zhao, K. Sakai, C. Blanch-Mercader, J. Prost, et al. 2023.
1058 “*Alcanivorax Borkumensis* Biofilms Enhance Oil Degradation by Interfacial Tubulation.”
1059 *Science (New York, N.Y.)* 381 (6659): 748–53.
- 1060 Price, Morgan N., Paramvir S. Dehal, and Adam P. Arkin. 2010. “FastTree 2--Approximately
1061 Maximum-Likelihood Trees for Large Alignments.” *PloS One* 5 (3): e9490.
- 1062 Rusch, Douglas B., Adam C. Martiny, Christopher L. Dupont, Aaron L. Halpern, and J. Craig
1063 Venter. 2010. “Characterization of *Prochlorococcus* Clades from Iron-Depleted Oceanic
1064 Regions.” *Proceedings of the National Academy of Sciences of the United States of
1065 America* 107 (37): 16184–89.
- 1066 Ruscheweyh, Hans-Joachim, Alessio Milanese, Lucas Paoli, Nicolai Karcher, Quentin Clayssen,
1067 Marisa Isabell Keller, Jakob Wirbel, et al. 2022. “Cultivation-Independent Genomes Greatly
1068 Expand Taxonomic-Profiling Capabilities of mOTUs across Various Environments.”
1069 *Microbiome* 10 (1): 212.
- 1070 Sabirova, Julia S., Tatyana N. Chernikova, Kenneth N. Timmis, and Peter N. Golyshin. 2008.
1071 “Niche-Specificity Factors of a Marine Oil-Degrading Bacterium *Alcanivorax Borkumensis*
1072 SK2.” *FEMS Microbiology Letters* 285 (1): 89–96.
- 1073 Salazar, Guillem, Lucas Paoli, Adriana Alberti, Jaime Huerta-Cepas, Hans-Joachim
1074 Ruscheweyh, Miguelangel Cuenca, Christopher M. Field, et al. 2019. “Gene Expression
1075 Changes and Community Turnover Differentially Shape the Global Ocean
1076 Metatranscriptome.” *Cell* 179 (5): 1068–83.e21.
- 1077 Sánchez, Pablo, Felipe H. Coutinho, Marta Sebastián, Massimo C. Pernice, Raquel Rodríguez-
1078 Martínez, Guillem Salazar, Francisco Miguel Cornejo-Castillo, et al. 2024. “Marine
1079 Picoplankton Metagenomes and MAGs from Eleven Vertical Profiles Obtained by the
1080 Malaspina Expedition.” *Scientific Data* 11 (1): 154.

- 1081 Sánchez, Pablo, Marta Sebastián, Massimo Pernice, Raquel Rodríguez-Martínez, Stephane
1082 Pesant, Susana Agustí, Takashi Gojobori, et al. 2023. "Marine Picoplankton Metagenomes
1083 from Eleven Vertical Profiles Obtained by the Malaspina Expedition in the Tropical and
1084 Subtropical Oceans." *bioRxiv*. <https://doi.org/10.1101/2023.02.06.526790>.
- 1085 Schattenhofer, Martha, Bernhard M. Fuchs, Rudolf Amann, Mikhail V. Zubkov, Glen A. Tarran,
1086 and Jakob Pernthaler. 2009. "Latitudinal Distribution of Prokaryotic Picoplankton
1087 Populations in the Atlantic Ocean." *Environmental Microbiology* 11 (8): 2078–93.
- 1088 Schmidt, Thomas S. B., Anthony Fullam, Pamela Ferretti, Askarbek Orakov, Oleksandr M.
1089 Maistrenko, Hans-Joachim Ruscheweyh, Ivica Letunic, et al. 2024. "SPIRE: A Searchable,
1090 Planetary-Scale microbiome REsource." *Nucleic Acids Research* 52 (D1): D777–83.
- 1091 Schroer, William Francis. 2023. "Metabolite Transport and Its Role in Marine Microbial
1092 Interactions." University of Georgia. <https://www.proquest.com/docview/2917419506>.
- 1093 Septer, Alecia N., and Karen L. Visick. 2024. "Lighting the Way: How the *Vibrio Fischeri* Model
1094 Microbe Reveals the Complexity of Earth's 'Simplest' Life Forms." *Journal of Bacteriology*
1095 206 (5): e0003524.
- 1096 Shaiber, Alon, Amy D. Willis, Tom O. Delmont, Simon Roux, Lin-Xing Chen, Abigail C. Schmid,
1097 Mahmoud Yousef, et al. 2020. "Functional and Genetic Markers of Niche Partitioning
1098 among Enigmatic Members of the Human Oral Microbiome." *Genome Biology* 21 (1): 292.
- 1099 Spang, Anja, Jimmy H. Saw, Steffen L. Jørgensen, Katarzyna Zaremba-Niedzwiedzka, Joran
1100 Martijn, Anders E. Lind, Roel van Eijk, Christa Schleper, Lionel Guy, and Thijs J. G. Ettema.
1101 2015. "Complex Archaea That Bridge the Gap between Prokaryotes and Eukaryotes."
1102 *Nature* 521 (7551): 173–79.
- 1103 Steinegger, Martin, and Johannes Söding. 2017. "MMseqs2 Enables Sensitive Protein
1104 Sequence Searching for the Analysis of Massive Data Sets." *Nature Biotechnology*,
1105 October. <https://doi.org/10.1038/nbt.3988>.
- 1106 Stepanauskas, Ramunas, Elizabeth A. Fergusson, Joseph Brown, Nicole J. Poulton, Ben
1107 Tupper, Jessica M. Labonté, Eric D. Becraft, et al. 2017. "Improved Genome Recovery and
1108 Integrated Cell-Size Analyses of Individual Uncultured Microbial Cells and Viral Particles."
1109 *Nature Communications* 8 (1): 84.
- 1110 Sul, Woo Jun, Thomas A. Oliver, Hugh W. Ducklow, Linda A. Amaral-Zettler, and Mitchell L.
1111 Sogin. 2013. "Marine Bacteria Exhibit a Bipolar Distribution." *Proceedings of the National
1112 Academy of Sciences of the United States of America* 110 (6): 2342–47.
- 1113 Sunagawa, Shinichi, Luis Pedro Coelho, Samuel Chaffron, Jens Roat Kultima, Karine Labadie,
1114 Guillem Salazar, Bardya Djahanschiri, et al. 2015. "Ocean Plankton. Structure and Function
1115 of the Global Ocean Microbiome." *Science* 348 (6237): 1261359.
- 1116 Thomas, Mridul K., Colin T. Kremer, Christopher A. Klausmeier, and Elena Litchman. 2012. "A
1117 Global Pattern of Thermal Adaptation in Marine Phytoplankton." *Science* 338 (6110): 1085–
1118 88.
- 1119 Tully, Benjamin J., Elaina D. Graham, and John F. Heidelberg. 2018. "The Reconstruction of
1120 2,631 Draft Metagenome-Assembled Genomes from the Global Oceans." *Scientific Data* 5
1121 (January):170203.
- 1122 Ustick, Lucas J., Alyse A. Larkin, and Adam C. Martiny. 2023. "Global Scale Phylogeography of
1123 Functional Traits and Microdiversity in *Prochlorococcus*." *The ISME Journal* 17 (10): 1671–
1124 79.
- 1125 Utter, Daniel R., Gary G. Borisy, A. Murat Eren, Colleen M. Cavanaugh, and Jessica L. Mark
1126 Welch. 2020. "Metapangenomics of the Oral Microbiome Provides Insights into Habitat
1127 Adaptation and Cultivar Diversity." *Genome Biology* 21 (1): 293.
- 1128 Watterson, William J., Melikhan Tanyeri, Andrea R. Watson, Candace M. Cham, Yue Shan,
1129 Eugene B. Chang, A. Murat Eren, and Savaş Tay. 2020. "Droplet-Based High-Throughput
1130 Cultivation for Accurate Screening of Antibiotic Resistant Gut Microbes." *eLife* 9 (June).
1131 <https://doi.org/10.7554/eLife.56998>.

- 1132 Weinheimer, Alaina R., and Frank O. Aylward. 2020. "A Distinct Lineage of Caudovirales That
1133 Encodes a Deeply Branching Multi-Subunit RNA Polymerase." *Nature Communications* 11
1134 (1): 4506.
- 1135 Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy McGowan, Romain
1136 François, Garrett Golemund, et al. 2019. "Welcome to the Tidyverse." *Journal of Open
1137 Source Software* 4 (43): 1686.
- 1138 Woodcroft, Ben J., Samuel T. N. Aroney, Rossen Zhao, Mitchell Cunningham, Joshua A. M.
1139 Mitchell, Linda Blackall, and Gene W. Tyson. 2024. "SingleM and Sandpiper: Robust
1140 Microbial Taxonomic Profiles from Metagenomic Data." *bioRxiv*.
1141 <https://doi.org/10.1101/2024.01.30.578060>.
- 1142 Wood, Derrick E., and Steven L. Salzberg. 2014. "Kraken: Ultrafast Metagenomic Sequence
1143 Classification Using Exact Alignments." *Genome Biology* 15 (3): R46.
- 1144 Woyke, Tanja, Devin F. R. Doud, and Frederik Schulz. 2017. "The Trajectory of Microbial
1145 Single-Cell Sequencing." *Nature Methods* 14 (11): 1045–54.
- 1146 Wu, Martin, and Jonathan A. Eisen. 2008. "A Simple, Fast, and Accurate Method of
1147 Phylogenomic Inference." *Genome Biology* 9 (10): R151.

1148