# New isolate genomes and global marine metagenomes resolve ecologically relevant units of SAR11

Kelle C. Freel[1], Sarah J. Tucker[1,2,3,4,5], Evan B. Freel[1], Stephen J. Giovannoni[6], A. Murat Eren,[3,4,5,7,8] & Michael S. Rappé[1*]

[1]Hawaiʻi Institute of Marine Biology, University of Hawaiʻi at Mānoa, Kāneʻohe, Hawaiʻi, United States; [2]Marine Biology Graduate Program, University of Hawaiʻi at Mānoa, Honolulu, Hawaiʻi, United States; [3]Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, MA 02543, USA; [4]Helmholtz Institute for Functional Marine Biodiversity, 26129 Oldenburg, Germany; [5]Alfred Wegener Institute, Helmholtz Centre for Polar and Marine Research, 27570 Bremerhaven, Germany; [6]Department of Microbiology, Oregon State University, 97331 Corvallis, OR, United States; [7]Institute for Chemistry and Biology of the Marine Environment, University of Oldenburg, 26129 Oldenburg, Germany; [8]Max Planck Institute for Marine Microbiology, 28359 Bremen, Germany

*Corresponding Author. Email: rappe@hawaii.edu

Running Title: SAR11 genomes from the tropical Pacific

1

# Abstract

The bacterial order *Pelagibacterales* (SAR11) is among the most abundant and widely distributed microbial lineages across the global surface ocean, where it forms an integral component of the marine carbon cycle. However, the limited availability of high-quality genomes has hampered comprehensive insights into the ecology and evolutionary history of this critical group. Here, we increase the number of complete SAR11 isolate genomes fourfold by describing 81 new SAR11 strains from seven distinct lineages isolated from coastal and offshore surface seawater of the tropical Pacific Ocean. We leveraged comprehensive phylogenomic insights afforded by these isolates to characterize 24 monophyletic, discrete ecotypes with unique spatiotemporal patterns of distribution across the global ocean, which we define as genera. Our data illustrate fine-scale differentiation in patterns of detection with ecologically-relevant gene content variation for some closely related genomes, demonstrating instances of ecological speciation within SAR11 genera. Our study provides unique insight into complex environmental SAR11 populations, and proposes an ecology-informed hierarchy to pave a path forward for the systematic nomenclature for this clade.

# Main

SAR11 marine bacteria are a genetically diverse, order-level lineage of heterotrophs within the *Alphaproteobacteria* known as the *Pelagibacterales* (Grote et al. 2012) that numerically dominate planktonic communities across the global ocean (Morris et al. 2002; Carlson et al. 2009; Eiler et al. 2009; Schattenhofer et al. 2009; Becker et al. 2019). Associations between the spatiotemporal distribution of operationally defined subclades and environmental variables suggest the presence of distinct ecotypes within SAR11 (Carlson et al. 2009; Eren et al. 2013a; Delmont et al. 2019; Tucker et al. 2021). Previous studies further support the functional differentiation of subclades (Grote et al. 2012; Thrash et al. 2014), even across short biogeographical distances (Tucker et al. 2024a). While limited in number, the available high-quality SAR11 genomes have demonstrated that this group is a remarkably cohesive genetic assemblage (Grote et al. 2012), making it an attractive model to study the capacity of a minimalist genome to reach stunning levels of success.

Since the first observation of SAR11 through environmental 16S rRNA gene fragments over three decades ago (Giovannoni et al. 1990), microbiology has benefited from a dramatic increase in microbial sequence data recovered directly from the environment, offering representative genomes for many difficult to cultivate microbial lineages (Hug et al. 2016). However, even the most comprehensive genome-resolved surveys of marine metagenomes have failed to yield high-quality SAR11 genomes (Paoli et al. 2022), resulting in limited insights into what constitutes ecologically meaningful units within this broad group. The extensive intra-clade diversity of SAR11 (Tsementzi et al. 2016; Kiefl et al. 2023) confounds the ability to reconstruct

3

environmental genomes from metagenomes (Delmont et al. 2018; Tully et al. 2018), which is why one of the most abundant microbial clades in marine systems suffers from poor representation in genome-resolved metagenomics surveys (Chang et al. 2024). Circumventing the need to assemble complex metagenomes first for genome recovery, single-cell sorting techniques have been much more effective in sampling environmental SAR11 populations through single-amplified genomes (SAGs). However, in an extensive effort to characterize surface ocean microbes, the estimated genome completion of SAGs that could be affiliated with SAR11 remained below 60% (Pachiadaki et al. 2019), a level that prevents robust phylogenomic insights. Such barriers have led to a reliance on isolate genomes to investigate the evolution of SAR11 populations (Vergin et al. 2007; Wilhelm et al. 2007; Thrash et al. 2011; Grote et al. 2012; Muñoz-Gómez et al. 2019), yet following this path has been impeded by another formidable challenge: the difficulty of cultivating SAR11 in the laboratory, even with genomic insights regarding its unique growth requirements (Tripp et al. 2008; Carini et al. 2013; Sun et al. 2016).

The first successful cultivation of SAR11 in 2002 resulted in the isolation of *Pelagibacter ubique* strain HTCC1062 (Rappé et al. 2002), followed by the publication of its complete genome (Giovannoni et al. 2005). Over the past two decades, additional isolate genomes have been few, with only 25 currently available. Despite their rarity, high-quality genomes from isolated strains not only shed light on SAR11 biology (Schwalbach et al. 2010; Sun et al. 2011; Carini et al. 2013) and the origins of this lineage within the *Alphaproteobacteria* (Thrash et al. 2011; Grote et al. 2012; Muñoz-Gómez et al. 2019), but also have made it possible to establish key concepts in biology such as genome streamlining (Schwalbach et al. 2010; Sun et al. 2011;

4

Grote et al. 2012; Viklund et al. 2012; Giovannoni et al. 2014; Giovannoni 2017) and investigate the evolutionary processes that shape protein evolution (Delmont et al. 2019; Kiefl et al. 2023).

Here we report 81 high-quality genomes from SAR11 strains, increasing the number available for SAR11 isolates by fourfold, and leverage this new collection to build a robust genome phylogeny for the order *Pelagibacterales*. By incorporating publicly available, high-quality single-cell genomes and surface ocean metagenomes from both a steep, nearshore to open-ocean local environmental gradient and elsewhere from around the globe, we reveal cohesive patterns of genomic and ecotypic diversification. We propose a framework through which to characterize and interpret genome heterogeneity at multiple stages along the evolutionary history of SAR11 marine bacteria, and establish a roadmap for future efforts to organize this globally abundant bacterial clade.

# Results

**Eighty-one high-quality genomes sequenced from 206 newly isolated SAR11 strains and co-cultures**

Three dilution-to-extinction culturing experiments using surface seawater collected from nearshore and adjacent offshore environments of Oʻahu, Hawaiʻi, in the tropical Pacific Ocean resulted in 916 isolates from 2,102 inoculated cultures (Table 1; Supplemental Fig. 1). Using a streamlined isolate-to-genome approach, we identified 206 cultures as either pure SAR11 strains or mixed cultures with at least 50% of the total reads matching a SAR11 strain via 16S rRNA gene amplicon sequencing (Supplemental Table 1), and sequenced draft genomes from 90. Manual curation resulted in 79 high-quality SAR11 isolate genomes. The genomes from two

100 strains (HIMB123 and HIMB109) isolated from a previous culture experiment were also added

101 (Brandon 2006), resulting in 81 new SAR11 genomes from isolates. The majority of these

102 (n=60) assembled into ten contigs or less, including 24 closed genomes and an additional 30

103 containing one to three contigs. They ranged from 1.00 to 1.54 Mbp in size and GC content of

104 28.5 to 30.7% (Supplemental Table 2). The median pairwise genome-wide average nucleotide

105 identity (gANI) value across all genomes was 81.8% and none of the 81 new isolate genomes

106 were identical. Having captured a genetically diverse array of SAR11 isolates, we used a

107 phylogenomic approach to characterize evolutionary relationships between these genomes and to

108 high-quality single-cell and isolate genomes previously retrieved from seawater.

109
110 **Table 1. Summary of high-throughput culturing (HTC) experiments.**
111

| Site | Inoculum source | Inoculum size (# of cells) | Cultures screened | Positive cultures | SAR11 genomes |
|---|---|---|---|---|---|
| SB | raw seawater | 5 | 576 | 339 | 53 |
| STO1 | raw seawater | 5 | 576 | 126 | 16 |
| STO1 | cryopreserved seawater | 5 | 480 | 142 | 9 |
| STO1 | cryopreserved seawater | 100 | 470 | 343 | 1 |

112

113 **A comprehensive genome phylogeny reveals a robust evolutionary backbone populated by**

114 **clusters of closely related genomes**

115    We first sought to resolve relationships between the strains isolated in this study and

116 other publicly available high-quality *Pelagibacterales* genomes to precisely establish where the

117 new genomes originate from within the broad spectrum of known SAR11 diversity. For this, we

118 created a database that, in addition to the 81 genomes presented here, included 25 public SAR11

119 isolate genomes, 8 of which were also isolated from off the windward coast of Oʻahu, Hawaiʻi,

120 and 375 SAR11 single-amplified genomes (SAGs) estimated to be ≥85% complete with a

6

121 redundancy <5% (Supplemental Table 3). We also included five additional SAR11 SAGs of

122 potentially unique evolutionary origin in this collection (Vergin et al. 2013; Thrash et al. 2014),

123 though we excluded genomes from putative SAR11 subgroups IV (Vergin et al. 2013) and V

124 (Thrash et al. 2011) due to their unlikely or, at a minimum, uncertain shared common ancestry

125 with SAR11 (Thrash et al. 2011; Viklund et al. 2013; Haro-Moreno et al. 2020; Muñoz-Gómez et

126 al. 2022). This resulted in a curated collection of 481 SAR11 genomes to assess the evolutionary

127 backbone for SAR11.

128        Previous studies investigating phylogenomic relationships within the

129 *Alphaproteobacteria* utilized a curated set of 200 single-copy core genes (SCGs) for this

130 bacterial class (Wang and Wu 2013; Muñoz-Gómez et al. 2019). We evaluated the presence of

131 these 200 SCGs across our genome dataset, and excluded genes missing in more than 90% of the

132 481 SAR11 genomes. This resulted in a SAR11-specific SCG set of 165 genes for downstream

133 phylogenomic analyses, referred to hereafter as the SAR11_165 core gene set (Supplemental
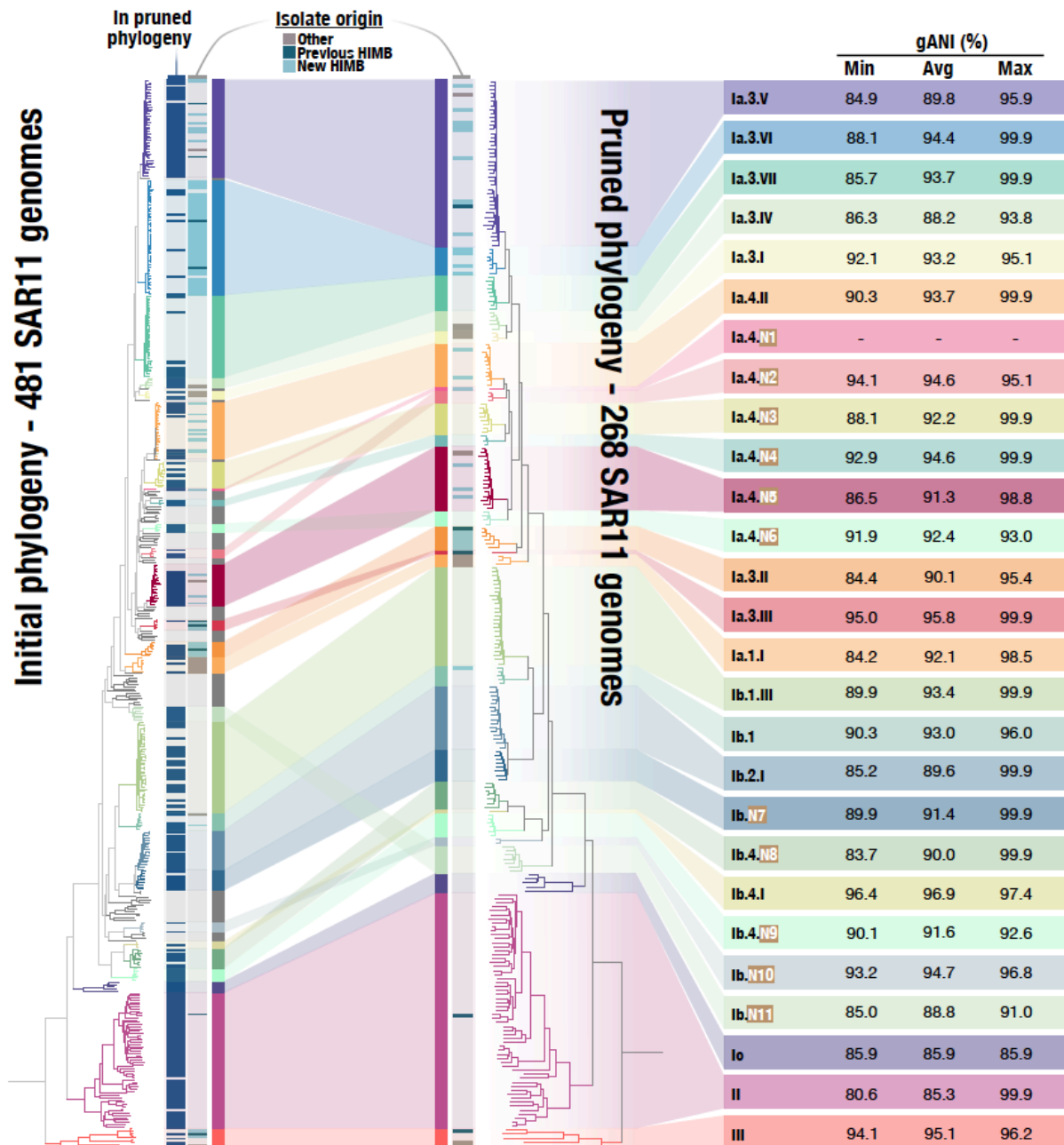
134 Table 4).

135        Our analysis of the 481 genomes using the SAR11_165 gene set revealed that the SAR11

136 clade consists of four robust, deeply-branching sublineages (Fig. 1; Supplementary Fig. 2). Three

137 of these branches were the previously characterized subclades Ic (Vergin et al. 2013), II (Suzuki

138 et al. 2001), and III (Morris et al. 2005), while the fourth was a combination of established

139 SAR11 subclades Ia and Ib (Suzuki et al. 2001), which did not form separate monophyletic

140 subclades in this comprehensive genomic dataset and robust phylogenetic analysis. If the SAR11

141 clade is assigned to the taxonomic level of a bacterial order, then these four lineages logically

142 resolve to the taxonomic level of families.

143    We further removed genomes from this initial tree in two steps. First we excluded SAGs

144 that did not fall into a 90% gANI cluster of at least three genomes to focus our analyses on

145 well-resolved regions of the tree. Second, we de-replicated the remaining genomes using a

146 conservative cutoff of 95% gANI to minimize subsequent competitive metagenomic recruitment

147 steps splitting reads among closely related genomes (Evans and Denef 2020). While the 95%

148 ANI cutoff is broadly recognized in contemporary microbiology as a threshold to identify

149 microbial species, it over-splits ecologically and evolutionarily cohesive units in SAR11 and

150 does not delineate species-like groups. We note that the reason behind our use of the 95% ANI in

151 this step of our analysis was solely to establish a technically robust workflow prior to

152 competitive read recruitment rather than a biologically meaningful partitioning of our genomes, a

153 challenge our study focuses on later.

154    We then turned our attention to the distal end of the phylogeny, which contained a large

155 number of well supported clusters of closely related genomes, particularly within the Ia/Ib

156 subgroup that contained 78 of the 81 new isolate genomes. A phylogeny of the resulting 268

157 genomes revealed 24 monophyletic clusters within the historical Ia/Ib subgroup that were

158 characterized by a range of gANI values from 84% to 96% (92.1 ± 2.94%; mean ± SD) (Fig. 1,

159 Supplemental Fig. 3). While a handful of these clusters were recognized previously, we defined

160 an additional 11 here (Fig 1; Supplemental Table 5). Twelve of the 24 clusters contained an

161 isolated representative, and eight contained at least one isolate from our study area in the tropical

162 Pacific.

163    In summary, our extensive phylogenomic analysis of SAR11 revealed 24 monophyletic

164 clusters within the historical Ia/Ib subgroup which included the majority of SAR11 SAGs and the

8

165 new and previously published isolate genomes. The non-uniform minimum gANI estimates

166 suggest that the application of sequence-based ANI thresholds to demarcate SAR11 diversity

167 may obscure important evolutionary signals. Hypothesized drivers of the maintenance and

168 partitioning of genomic diversity in SAR11 include niche-based processes, where genetically

169 cohesive clusters also display ecological homogeneity and the underlying genetic diversity is

170 maintained by similar forces of selection, recombination, and drift. To understand the potential

171 eco-evolutionary forces that shape SAR11 diversification, we turned to metagenomic read

172 recruitment analysis to recover biogeographical distribution patterns for our genomes across the

173 globe.

**Fig 1. Comprehensive phylogenies of the *Pelagibacterales*.** A comparison between an exhaustive phylogeny (left panel) with 481 SAR11 genomes (106 isolates and 375 SAGs) and a pruned phylogeny (right panel) with 268 genomes (50 isolates and 218 SAGs), based on a curated SAR11-specific set of 165 genes. Genomes included in the pruned phylogeny are indicated with a dark blue bar in the left panel, and the origin of isolate genomes is indicated for both phylogenies.

**Global read recruitment from the surface ocean reveals broadly congruent phylogenetic and ecotypic diversification across SAR11**

Our competitive metagenomic read recruitment assessed the distribution of the 268 SAR11 genomes around the globe and relied upon 950 publicly-available marine metagenomes, as well as metagenomes from the Kāneʻohe Bay Time-series (KByT), the location of isolation for the 81 new and 9 of the 25 existing isolate genomes (Supplemental Table 6; Supplemental Table 7). These data enabled us to investigate whether cohesive genomic and ecological groups, or ecotypes, could be discerned by combining SAR11 phylogeny and biogeography.

Our first priority was to establish whether genome clusters within a given SAR11 clade showed cohesive read recruitment profiles across metagenomes, or, in other words, whether the ecological patterns revealed by a single genome were similar to all genomes within the group to which it belonged. Detection values for multiple genomes within a genome cluster showed a high degree of cohesion (Fig. 2; Supplemental Table 8; Supplemental Table 9). For example, representatives from Ia.3.IV, Ia.3.I, Ia.4_II, Ia.4.N2, Ia.4.N5, Ib.1.III, and Ib.4.N9 are particularly consistent within the genome clusters (Fig. 2). Consistent overlap between SAGs and isolate genomes within the same clade demonstrated that both genome types accurately reflect distribution patterns for closely related populations as inferred by phylogeny (Fig. 2). A non-metric multidimensional scaling (NMDS) analysis of the overall detection patterns of genomes across metagenomes consistently grouped genomes within a given clade more closely compared to those that belonged to other genome clusters (Supplemental Fig. 4), further supporting a high degree of intra-clade ecological cohesion.

11

203    Our second priority was to establish insights into whether SAR11 genome clusters

204 differed in their biogeographical patterns, and whether genome clusters identified SAR11

205 populations of distinct ecology. Hierarchical clustering of metagenomes based on SAR11

206 detection patterns revealed four groups: metagenomes that originated from (1) low-latitude

207 samples, (2) high-latitude samples with low SAR11 diversity, (3) low-latitude samples with high

208 SAR11 diversity, as well as (4) samples from coastal Kāneʻohe Bay (Fig. 2). Many SAR11

209 genome clusters were indeed differentially distributed across these metagenome groups. For

210 example, Ia.4.II and Ia.4.N5 were only consistently found in groups 3 and 4, while Ia.3.IV was

211 found across group 1 and only in select sites in groups 3 and 4 (Fig 2). However, in multiple

212 cases, the environmental detection patterns of different phylogenomic genome clusters

213 overlapped; while there was some degree of inter-clade ecological differentiation, distinct

214 SAR11 genome clusters frequently co-occurred (Fig. 2, Supplemental Fig. 5). This observation

215 suggests that patterns of distribution alone cannot discern the boundaries of cohesive ecological

216 units within SAR11, a task that evidently requires the integration of biogeographical patterns

217 through metagenomic read recruitment with ancestral relationships among genomes though

218 phylogenomics.

219    Finally, we used the read recruitment analysis to assign ecological patterns to specific

220 SAR11 genome clusters. While multiple broad patterns were clear from the pairing of the

221 phylogenomic relationships and read recruitment data, we focused our investigation on whether

222 the genome clusters within the SAR11 Ia/Ib lineage that appeared to be confined to the coastal

223 end of the KByT environmental gradient (Ia.3.VI, Ia.3.II, and Ia.3.III; Supplemental Fig. 1; also

224 see (Tucker et al. 2024a)) were similarly constrained to coastal areas globally. Indeed, two of the

12

225 three genome clusters, Ia.3.II and Ia.3.III, were detected almost exclusively in metagenomes

226 sourced from coastal environments (e.g., KByT, the north coast of Panama, the Chesapeake Bay,

227 and the Atlantic coast of Portugal). Interestingly, while the clade Ia.3.VI was restricted to

228 nearshore metagenomes across KByT, it was well-detected in both coastal and offshore

229 environments in other oceanic regions (Fig. 2). Genome clusters Ia.3.II and Ia.3.III did not

230 include any SAGs and were only composed of isolates from coastal Kāneʻohe Bay. Yet, we could

231 detect them in other oceans, which confirms their global relevance as representatives of SAR11

232 populations adapted to coastal ecosystems.

233    Through the combination of global metagenomic read recruitment and phylogenomics,

234 we show that SAR11 genome clusters contain genomes with a high degree of intra-clade

235 ecological cohesion. These genome clusters were often distinguished by their ecological

236 distributions and demonstrated notable inter-clade ecological differentiation. Finally, we applied

237 this framework to understand how SAR11 genetic and ecological diversity partitions among

238 ocean biomes, in particular coastal ocean and open ocean environments.

239    The integrated ecological and evolutionary framework here is supported by high-quality

240 genomes that span the known diversity of the *Pelagibacterales*, providing a critical opportunity

241 to discern distinct ecologically meaningful genome clusters within SAR11. We show that the 24

242 distinct genome clusters represent groups sharing cohesive ecological patterns and evolutionary

243 relationships, not at the finest tips of the phylogenomic tree, but at relatively deeper branches

244 that encompass gANI values ranging between 84% and 96%. This suggests it is unlikely that

245 these genome clusters represent SAR11 diversity at the level of 'species'. This conclusion is

246 further supported by our companion work (Tucker et al. 2024a), which reveals systematic

13

247 differences in the metabolic potential of SAR11 genome clusters that likely support distinct

248 ecological distributions in immediately adjacent coastal and open ocean surface seawater with

249 habitat-specific metabolic genes that are under higher selective forces. With the combined

250 evidence presented here and in the work by Tucker et al. (2024a) that unite SAR11 diversity into

251 distinct genome clusters with ecotype properties supported by SAR11 phylogenomics, ecology,

252 metabolic potential, as well as population genetics, we argue that the most conceivable

253 taxonomic rank at which SAR11 genome clusters can be described in a conventional framework

254 emerges as the 'genus' level.

255     This genus-level designation is ideal as it encompasses a degree of diversity previously

256 designated by SAR11 subgroups and has the flexibility to account for subtle variation in ecology

257 recognized between closely related genomes. We identified the highest quality genome

258 representatives (electing for isolates when possible) to assign as type genomes for each genus

259 (Fig. 4), which establishes a roadmap to rationally designate new genera as they are identified in
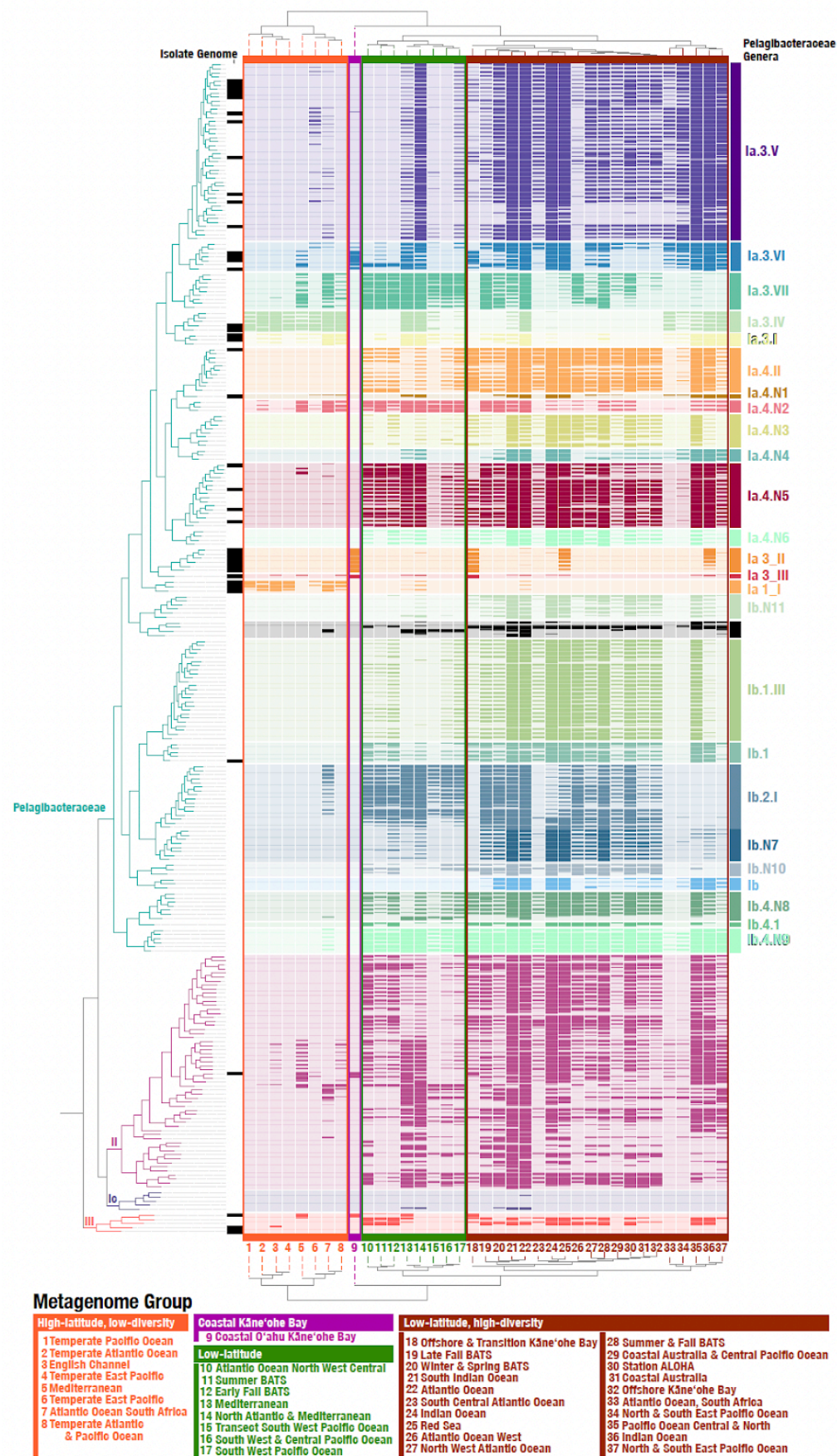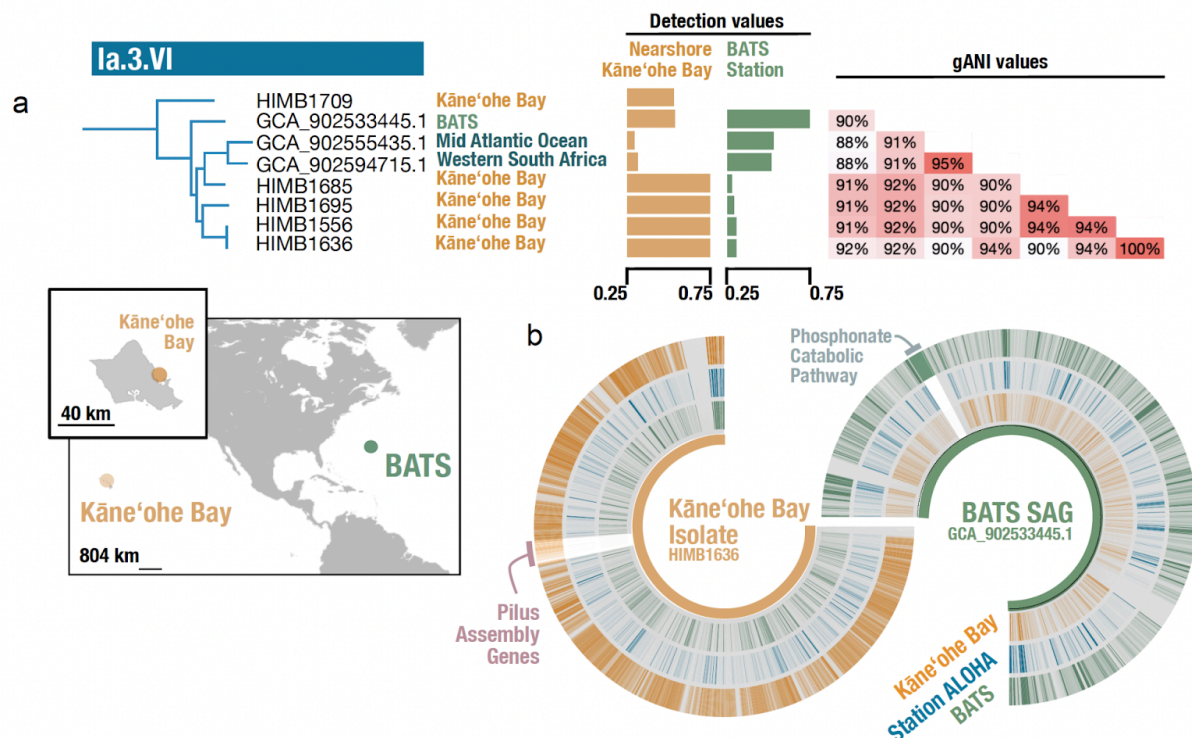
260 the future.

261

263 **Fig. 2: Global metagenome read recruitment to 268 *Pelagibacterales* genomes.** A clustering
264 analysis reveals that the distribution of metagenomes from the same geographic location have
265 characteristic patterns of detection. Detection values from 0.25 to 0.75 are shown.
266

267 **Evidence for ecological speciation within closely related genome clusters**

268       Despite broad ecological cohesion within what we have designated as *Pelagibacterales*

269 genera, some notable differences highlight underlying complexities in defining the finest scales

270 of divergence. The Ia.3.VI genus includes genomes from strains of Kāneʻohe Bay origin as well

271 as SAGs from other regions of the global ocean and encompasses significant genomic diversity

272 (minimum gANI 88%) and phylogenomic structure (Fig. 3a). Through read recruitment, we

273 observed notable differences in detection patterns of genomes across metagenomic samples.

274 Isolate genomes from the bay harbored the highest detection values of the Ia.3.VI genus from

275 metagenomes in the bay, while a SAG from the BATS site in the Atlantic Ocean

276 (GCA_902533445.1) had the highest detection values at the BATS site (Fig. 3a), particularly in

277 the summer and fall (Fig. 2).

278

279

16

280

**Fig. 3: Fine-scale ecological speciation between closely related SAR11 genomes. (a)** Detailed view of the Ia.3.VI genus including evolutionary relationships, locations of genome origin, detection values from select locations, and within-genus gANI values. The geographic origins of two of the closely related genomes that have distinct detection patterns include Kāneʻohe Bay in the Pacific Ocean and the Bermuda Atlantic Time-series Study (BATS) in the Atlantic. **(b)** Coverage values of isolate HIMB1636 and SAG GCA_902533445.1 of metagenomes from nearshore Kāneʻohe Bay, Station ALOHA in the North Pacific Subtropical Gyre, and BATS highlighting the differential detection of genes for type IV pilus assembly and the phosphonate catabolic pathway.

290

Given the underlying genomic diversification between isolate HIMB1636 and BATS

SAG GCA_902533445.1 and their distinct biogeographical distributions that peak in each of

their respective source locations, we next surveyed the genomes for potentially unique metabolic

capabilities. By inspecting the coverage of isolate HIMB1636 and BATS SAG

GCA_902533445.1 using metagenomes from Kāneʻohe Bay and the BATS site (Supplemental

Table 10), we found one genomic region of SAG GCA_902533445.1 that had particularly high

17

297 coverage at BATS compared to the KByT samples and included 29 genes encoding the uptake

298 (*phnCDE*) and catabolism of phosphonates via the C-P lyase pathway (*phnGHIJKLM*) (Fig. 3d)

299 (Villarreal-Chiu et al. 2012). The *phnJ* phylogeny did not reflect the phylogenomic relationships

300 among genomes (Supplemental Fig. 6), and the entire pathway was located on a genomic island

301 similar to the marine bacterium HIMB59 (Molina-Pardines et al. 2023). The C-P lyase pathway

302 is known to be enriched in phosphate-depleted systems of the Atlantic Ocean (Sosa et al. 2019;

303 Acker et al. 2022), so the presence of the C-P lyase catabolic genes in a genome sourced from

304 BATS, but missing from a closely related genome sourced from more phosphate-replete

305 environments of Kāneʻohe Bay in the Pacific, suggests these genes provide an advantage in

306 phosphate depleted systems and that the BATS SAG GCA_902533445.1 may be locally-adapted

307 to these environments.

308　　　　While the HIMB1636 genome lacked the C-P lyase pathway, it contained a unique

309 genomic region with particularly high coverage that was not found in SAG GCA_902533445.1,

310 and encoded genes for type IV pilus assembly. The role of type IV pilus assemblies in SAR11 is

311 unclear (Zhao et al., 2017), although in other organisms it has been associated with an array of

312 functions including DNA uptake, twitching motility, and aggregation into microcolonies (Craig

313 and Li 2008). The presence of the type IV pilus assembly genes in the Ia.3.VI genome sourced

314 from the nitrogen-limited Pacific Ocean, but not in genomes from relatively more

315 nitrogen-replete waters of BATS, along with evidence that the *Pelagibacteraceae* can utilize

316 purine nucleosides and purine-derivatives for nitrogen (Braakman et al. 2024; Tucker et al.

317 2024a), suggests that the presence of a type IV pilus may be advantageous for DNA uptake in

318 nitrogen-poor environments and that HIMB1636 may be locally-adapted. Contrary to the
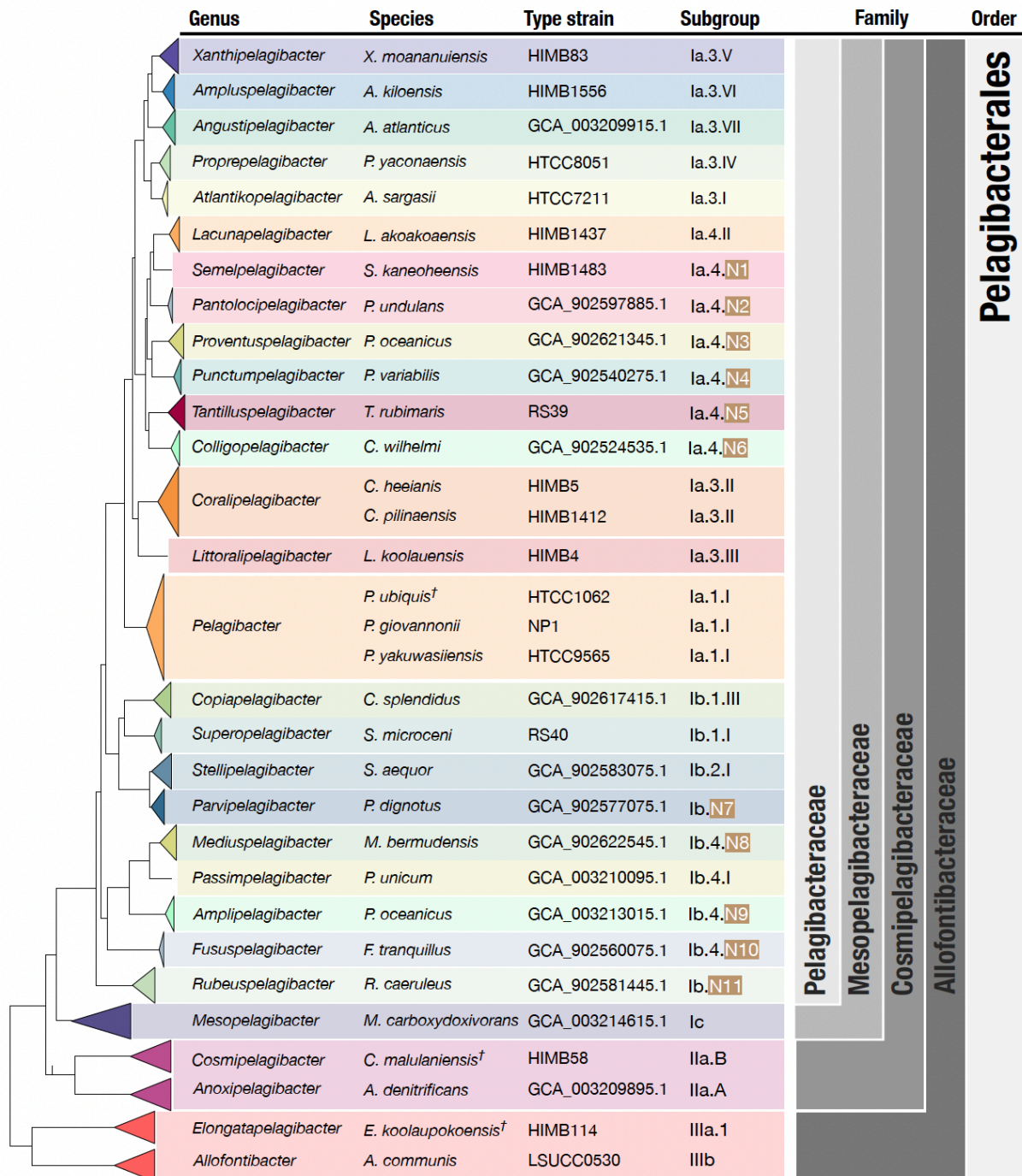
18

319 hypothesis that genera recombine at a rate sufficient to limit the ecotypic diversification of

320 closely related genomes (Zhao et al. 2024), our read mapping instead shows that the HIMB1636

321 and SAG GCA_902533445.1 genomes within cluster Ia.3.VI have sufficiently diverged at the

322 nucleotide level to reveal clear biogeographic divergence, and that they possess sets of genes that

323 reside in hypervariable genomic regions that are clearly associated with the differences in

324 abundance.

325     We examined gANI estimates, phylogenetic branching, environmental distributions, and

326 ecologically-relevant gene content to support the characterization of ecological diversification at

327 the finest tips of the tree, a process that we theorize to represent speciation. This underscores the

328 complexity of SAR11 ecology, highlights the need to include a diversity of representative

329 genomes within even closely related genera for environmental genomics studies, and indicates

330 that continued efforts to sample SAR11 globally are key to understanding the distribution of this

331 ubiquitous clade.

332

333 **Proposed *Pelagibacterales* classification and nomenclature**

334     We leveraged the robust genome phylogeny, gANI metrics, and read recruitment to

335 establish a rational classification and nomenclature system for the *Pelagibacterales* bacterial

336 order *Pelagibacterales*. To provide a framework and vocabulary to discuss groups of SAR11 in a

337 meaningful context, we first defined four family-level monophyletic groups as the

338 *Pelagibacteraceae* (historical subgroups Ia and Ib), *Cosmipelagibacteraceae* (historical

339 subgroup II), *Allofontibacteraceae* (historical subgroup III), and the *Mesopelagibacteraceae*

340 (historical subgroup Ic) (Fig. 4). We focused our efforts primarily on classification within the

19

341 *Pelagibacteraceae* where the majority of cultured isolates originate. Within the

342 *Pelagibacteraceae*, we used phylogenomics and ecological data to characterize 24 genera that

343 represent cohesive genetic and ecological clades, and designate type species for each

344 (Supplemental Table 11). The primary aim of these efforts is to ensure that the taxonomic

345 hierarchy for SAR11 provides a useful and tractable reflection of the ecology and genetic

346 diversity within this globally distributed group, and establishes a rational system that future

347 efforts can build upon.

| Genus | Species | Type strain | Subgroup | Family | Order |
|---|---|---|---|---|---|
| Xanthipelagibacter | X. moananuiensis | HIMB83 | Ia.3.V | | |
| Ampluspelagibacter | A. kiloensis | HIMB1556 | Ia.3.VI | | |
| Angustipelagibacter | A. atlanticus | GCA_003209915.1 | Ia.3.VII | | |
| Proprepelagibacter | P. yaconaensis | HTCC8051 | Ia.3.IV | | |
| Atlantikopelagibacter | A. sargasii | HTCC7211 | Ia.3.I | | |
| Lacunapelagibacter | L. akoakoaensis | HIMB1437 | Ia.4.II | | |
| Semelpelagibacter | S. kaneoheensis | HIMB1483 | Ia.4.N1 | | |
| Pantolocipelagibacter | P. undulans | GCA_902597885.1 | Ia.4.N2 | | |
| Proventuspelagibacter | P. oceanicus | GCA_902621345.1 | Ia.4.N3 | | |
| Punctumpelagibacter | P. variabilis | GCA_902540275.1 | Ia.4.N4 | | |
| Tantilluspelagibacter | T. rubimaris | RS39 | Ia.4.N5 | | |
| Colligopelagibacter | C. wilhelmi | GCA_902524535.1 | Ia.4.N6 | | |
| Coralipelagibacter | C. heeianis | HIMB5 | Ia.3.II | | |
| Coralipelagibacter | C. pilinaensis | HIMB1412 | Ia.3.II | | |
| Littoralipelagibacter | L. koolauensis | HIMB4 | Ia.3.III | Pelagibacteraceae | Pelagibacterales |
| Pelagibacter | P. ubiquis[†] | HTCC1062 | Ia.1.I | | |
| Pelagibacter | P. giovannonii | NP1 | Ia.1.I | | |
| Pelagibacter | P. yakuwasiiensis | HTCC9565 | Ia.1.I | | |
| Copiapelagibacter | C. splendidus | GCA_902617415.1 | Ib.1.III | | |
| Superopelagibacter | S. microceni | RS40 | Ib.1.I | | |
| Stellipelagibacter | S. aequor | GCA_902583075.1 | Ib.2.I | | |
| Parvipelagibacter | P. dignotus | GCA_902577075.1 | Ib.N7 | | |
| Mediuspelagibacter | M. bermudensis | GCA_902622545.1 | Ib.4.N8 | | |
| Passimpelagibacter | P. unicum | GCA_003210095.1 | Ib.4.I | | |
| Amplipelagibacter | P. oceanicus | GCA_003213015.1 | Ib.4.N9 | | |
| Fususpelagibacter | F. tranquillus | GCA_902560075.1 | Ib.4.N10 | | |
| Rubeuspelagibacter | R. caeruleus | GCA_902581445.1 | Ib.N11 | | |
| Mesopelagibacter | M. carboxydoxivorans | GCA_003214615.1 | Ic | Mesopelagibacteraceae | |
| Cosmipelagibacter | C. malulaniensis[†] | HIMB58 | IIa.B | Cosmipelagibacteraceae | |
| Anoxipelagibacter | A. denitrificans | GCA_003209895.1 | IIa.A | | |
| Elongatapelagibacter | E. koolaupokoensis[†] | HIMB114 | IIIa.1 | Allofontibacteraceae | |
| Allofontibacter | A. communis | LSUCC0530 | IIIb | | |

348

**Figure 4. A proposed taxonomic framework for the SAR11 order _Pelagibacterales_.** This schematic SAR11 phylogeny unites proposed genus and species names, proposed type strains, and historical reference labels.

21

# Discussion

By integrating high-throughput cultivation experiments with publicly available genomes and metagenomes, our study provides key insights into a long-standing question: to what extent, and at what hierarchical levels, can the genomic and ecological diversity of SAR11 be partitioned into cohesive units? Through comprehensive phylogenomic analyses paired with global metagenomic read recruitment surveys, we reveal ecotypic differentiation at both relatively shallow, species-level and deeper, genus-level diversity within SAR11. This robust eco-evolutionary framework, which unifies independent yet complementary approaches to genomic diversity and biogeography, resolves the order *Pelagibacterales* into four families and the family *Pelagibacteraceae* into 24 genera, establishing a much-needed taxonomic framework that delineates SAR11 diversity into tractable units and provides a foundation for future investigations.

A tight relationship between the phylogeny and ecology of SAR11 has long been suggested (Brown et al. 2012; Vergin et al. 2013); however, the ability to associate specific SAR11 clades with distinct ecological patterns and explain forces that maintain SAR11 diversity has remained elusive. Focusing on sequence-discrete groups within deep ocean SAR11 lineages, a recent study concluded that recombination, rather than ecological speciation, was likely the major driver of species-level cohesion (Zhao et al., 2024). While this observation may explain forces that maintain species-level cohesion for some populations in this group, our study shows that the global sampling of environmental populations through metagenomes consistently supports ecological delineations that are congruent with phylogenomic clustering patterns,

22

373 pointing towards ecotypic differentiation as the pervasive driver of the evolution within the

374 *Pelagibacterales*. Interestingly, SAR11 genera that showed similar biogeographical distribution

375 patterns in our analysis tended to occupy distant parts of the tree. This observation suggests an

376 inverse correlation between the genetic similarity among SAR11 populations and their

377 co-occurrence, a trend known as phylogenetic overdispersion. Phylogenetic overdispersion has

378 been observed across the tree of life (Davies 2006) and is driven by forces of competitive

379 exclusion, an overarching ecological phenomenon that limits the co-occurrence of ecologically

380 similar, closely related organisms. Future analyses that aim to resolve specific genetic

381 determinants of competitive exclusion or co-existence may benefit from geographically

382 constrained time-series data, as these patterns are likely not immediately attainable from global

383 yet spatiotemporally sparse metagenomes.

384       The practical need of microbiologists to find reasonable cutoffs to demarcate species

385 boundaries from genomic data alone and the nature of SAR11 evolution do not align seamlessly.

386 Through the analysis of genomes, a large number of anecdotal observations support 95% ANI as

387 a reasonable means to resolve archaeal and bacterial species (Jain et al. 2018; Olm et al. 2020).

388 However, SAR11 serves as a reminder that practical solutions do not necessarily apply to all

389 microbial clades (Delmont et al. 2019; López-Pérez et al. 2020). One of the implications of the

390 efforts to standardize the tree of life based on principles that work only for the majority of

391 microbial taxa is the conflation of all SAR11 genomes into two genera in the taxonomic

392 framework derived from genomes available on GTDB based on RED scores (Parks et al. 2022).

393 Indeed, while the ecologically relevant units of SAR11 described in our study are in agreement

394 with functional, evolutionary, and ecological observations, they are in disagreement with the

23

395 contemporary summaries of this clade based on RED- or ANI-based demarcations. The ways in

396 which evolutionary relationships between distinct clades of life intersect with taxonomic

397 classification systems will unlikely be resolved in a manner that satisfies everyone in

398 microbiology (Waite et al. 2020; Sanford et al. 2021). In this juncture, we believe that a stronger

399 motivation to understand the biological drivers that render SAR11 incompatible with our best

400 practical approaches will bring us closer to a unified solution to partition microbial diversity into

401 meaningful units, rather than casting SAR11, one of the most numerous microbial clades on our

402 planet, as a mere outlier.

403      Insights into the eco-evolutionary processes that shape SAR11 diversification in our

404 study rely heavily on the contribution of 81 new isolate genomes that represent abundant and

405 ecologically-relevant SAR11 populations across the coastal and global ocean. The

406 ecology-informed hierarchical organization of these genomes enabled us to propose SAR11

407 genera with formal names here, and investigate the likely functional determinants of ecological

408 diversification across the *Pelagibacteraceae* in our companion work (Tucker et al. 2024a). While

409 deeper understanding of the physiological, metabolic, and genetic factors that shape SAR11

410 biology will require controlled experimentation of isolated strains in the laboratory, our study

411 organizes the eco-evolutionary characteristics of known SAR11 diversity and provides a

412 roadmap for future efforts aimed to organize and understand the ubiquitous SAR11 populations

413 inhabiting the global ocean.

414

# Methods

**High-throughput culturing from surface seawater within and adjacent to Kāneʻohe Bay, Oʻahu**

Growth medium was prepared as previously described (Monaghan et al., 2020). Briefly, 20 L of seawater was collected first on 8 July 2017 and again on 20 September 2017 from a depth of 2 meters at station SR4 (N 21º 27.699', W 157º 47.010') in acid-washed polycarbonate bottles (Supplemental Fig. 1). The seawater was then filtered, autoclaved, and sparged as previously described (Monaghan et al. 2020). After processing, the sterile seawater was stored at 4ºC until use.

Two 4 L seawater samples to be used as inoculum were collected on 26 July 2017 in acid-washed polycarbonate bottles from 2 meters from stations SB (N 21° 26.181', W 157° 46.642) and STO1 (N 21° 28.974, W 157° 45.978') (Supplemental Fig. 1) and immediately returned to the laboratory for further processing. All of the Kāneʻohe Bay Time series sampling sites were previously classified as 'nearshore', 'transition', or 'offshore', with SB and STO1 representing nearshore and offshore sites, respectively (Tucker et al. 2021). Subsamples of the raw seawater were processed as described previously (Monaghan et al. 2020). Briefly, aliquots were taken for cryopreservation in a final concentration of 10% v/v glycerol and fixed with paraformaldehyde for the enumeration of planktonic microorganisms via flow cytometry. Additionally, 0.96 L from station SB and 1.30 L from station STO1 were filtered through a 25 mm diameter, 0.1 µm pore-sized polyethersulfone membrane (Supor-100; Pall Gelman Inc., Ann

25

435 Arbor, MI), which was then submerged in 500 µL DNA lysis buffer and stored at -80ºC until

436 DNA extraction.

437       Subsamples of raw seawater from SB and STO1 were enumerated using microscopy,

438 diluted to 2.5 cells mL$^{-1}$, and plated in 2 mL volumes into a total of 1,152 wells (576 wells per

439 site) of custom-fabricated 96-well Teflon microtiter plates. This experiment is referred to here as

440 HTC2017. Plates were then sealed with breathable polypropylene microplate adhesive film and

441 incubated in the dark at 27ºC. Plates were monitored for cellular growth at 3.5 and 8 weeks using

442 flow cytometry as previously described (Tripp et al. 2008; Monaghan et al. 2020). Wells with

443 positive growth (greater than $10^4$ cells mL) after 24 or 57 days of incubation were further

444 sub-cultured by transferring approximately 1 mL into 20 mL of sterile seawater media amended

445 as previously described (Monaghan et al. 2020) with 400 µM $(NH_4)_2SO_4$, 400 µM $NH_4Cl$, 50 µM

446 $NaH_2PO_4$, 1 µM glycine, 1 µM methionine, 50 µM pyruvate, 800 nM niacin (B3), 425 nM

447 pantothenic acid (B5), 500 nM pyridoxine (B6), 4 nM biotin (B7), 4 nM folic acid (B9), 6 µM

448 myo-inositol, 60 nM 4-aminobenzoic acid, and 6 µM thiamine hydrochloride (B1). These

449 subcultures were then incubated at 27ºC in the dark for an additional 33 days and then all

450 samples were processed and cataloged.

451       Cultures checked at 33 days that yielded positive growth ($>10^4$ cells ml$^{-1}$) were

452 cryopreserved in duplicate (2 x 500 µL culture and a final concentration of 10% v/v glycerol).

453 Each well with positive growth was assigned an HIMB culture ID and cells from the

454 approximately 18 mL remaining volume of each culture were collected by filtration through a 13

455 mm diameter, 0.03 µm pore-sized polyethersulfone membrane (Sterlitech, Kent, WA, USA),

456 which was then submerged in 250 µL DNA lysis buffer and stored at -80ºC until DNA

457 extraction. The lysis buffer was prepared by adding the following to MilliQ water: 8 mL 1M Tris

458 HCl (pH 8.0), 1.6 ml 0.5M EDTA (pH 8.0), and 4.8 g Triton X, for a final volume of 400 mL,

459 which was then filter sterilized, with lysozyme added to aliquots immediately before use (at a

460 final concentration of 20 mg ml$^{-1}$).

461       An additional experiment was performed using cryopreserved samples of seawater

462 collected on July 26, 2017, and described previously (Monaghan et al. 2020). Briefly, the

463 cryopreserved sample was enumerated and then diluted to two cell concentrations (2.5 and 52.5

464 cells ml$^{-1}$), and used to plate 480 and 470 2-ml dilution cultures, respectively. This experiment is

465 referred to as HTC2018. Growth was monitored at 2, 3, and 5 weeks after inoculation with

466 positive growth ($>10^4$ cells ml$^{-1}$) from the 2.5 cells ml$^{-1}$ cultures subcultured into 20 ml of sterile

467 seawater growth medium and monitored for growth for up to 10 weeks at 27ºC in the dark.

468 Subcultures were then cryopreserved and cells collected for DNA sequencing as described

469 above. One well from the 52.5 cells ml$^{-1}$ inoculation was directly collected for DNA sequencing

470 without subculturing (Monaghan et al. 2020).

471 **DNA extraction and 16S rRNA gene amplicon sequencing**

472       Genomic DNA (gDNA) from all filtered cultures as well as environmental DNA (eDNA)

473 from STO1 and SB was extracted using the Qiagen DNeasy Blood and Tissue Kit with modified

474 manufacturer's instructions for bacterial cells (Qiagen, Germantown, Maryland, USA). The

475 modifications included the addition of an initial freeze-thaw step (3 cycles of 10 minutes at 65ºC

476 followed by 10 minutes at -80ºC), the addition of 35 μL Proteinase K and 278 μL buffer AL at

477 the appropriate pretreatment step, and finally when eluted the same 200 μL volume was passed

478 through the membrane three times.

27

479     For the initial identification of all cultures, gDNA was used as template for the

480 polymerase chain reaction (PCR) amplification (Bio Rad C1000 Touch, Bio Rad, Hercules, CA,

481 USA) using barcoded 515F and 926R primers targeting the V4 region of the SSU rRNA gene

482 (Parada et al., 2016) in a reaction volume of 25 µL composed of: 2 µL gDNA, 0.5 µL each

483 forward and reverse primer, 10 µL 5PRIME HotMasterMix (Quantabio, Beverly, MA, USA),

484 and 12 µL of molecular grade $H_2O$ (Monaghan et al. 2020). The reaction was as follows: an

485 initial denaturation step of 3 min at 94ºC, 40 cycles of 45 sec at 94ºC followed by 1 min at 50ºC

486 and 1.5 min at 72ºC, with a final extension of 10 min at 72ºC. The PCR products were prepared

487 for sequencing as previously described (Monaghan et al. 2020) and sequenced on a MiSeq

488 platform by the Oregon State University Center for Genome Research and Biocomputing.

489 **16S rRNA gene sequence analysis**

490     Amplicon sequence data were processed as previously described (Monaghan et al. 2020).

491 Briefly, the data was imported into QIIME2 v2019.4.0, and demultiplexed before being assessed

492 for sequence quality and merged. DADA2 (Callahan et al. 2016) was then used for quality

493 control. Taxonomy was assigned to all reads using a Naïve Bayes classifier trained on the Silva

494 rRNA v132 database (Quast et al. 2013). Cultures were first classified as defined previously

495 (Monaghan et al. 2020), with "monocultures" consisting of more than 90% of reads from a single

496 amplicon sequence variant (ASV), "mixed cultures" with an ASV that was between 50% and

497 90% of the reads, and finally cultures with no dominant members. Any samples with less than

498 1,000 reads were not included in further analyses. We aimed to sequence all strains that included

499 monocultures and mixed cultures of SAR11.

28

**Genome sequencing**

To prepare samples of interest for whole genome sequencing, all extractions with gDNA concentrations above 0.06 ng µL$^{-1}$, a total of 10 µL were aliquoted for sequencing. For samples with concentrations below 0.06 ng µL$^{-1}$, the remaining extraction volume (approximately 175 to 185 µL) was concentrated using a SpeedVac (ThermoFisher) to approximately 30 µl and was re-quantified (Qubit 2.0, Invitrogen). From the concentrated samples with a minimum of 0.06 ng µL$^{-1}$, 10 µL was aliquoted for sequencing. Samples for sequencing were prepared using a Nextera library kit and sequenced on the NextSeq500 platform via a 150 bp paired-end run.

Genomes for previously cultured strains HIMB109 and HIMB123 (Brandon 2006) were sequenced by the Joint Genome Institute. Multiple methods were used to sequence these two strains, including directly using 200 µL of cell culture for library prep as well as using multiple volumes (5, 10, or 20 µL) of culture for multiple displacement analysis (MDA) prior to library preparation. The genomes were evaluated based on completeness, length, number of reads, and total contigs post assembly using SPAdes (Bankevich et al. 2012). An additional assembly using all reads generated from various sequencing attempts per genome was also constructed using the same assembly method, the highest quality genomes based on the metrics above were manually curated and used for additional analyses.

**Genome assembly and assessment**

Short reads were trimmed with Trim Galore! (https://github.com/FelixKrueger/TrimGalore) and assembled using Unicycler (Wick et al. 2017), which acts as a SPAdes (Bankevich et al. 2012) optimizer with Illumina short read data. Once assembled, reference indexes were built, and read mapping was performed using Bowtie2 with

29

522 default parameters (Langmead and Salzberg 2012). SAMtools (Li et al. 2009) was used to

523 convert the SAM file to a sorted and indexed BAM file. These initial assemblies and BAM files

524 were used to visualize genomes in anvi'o to check for possible contamination (Eren et al. 2015,

525 2021). For genomes with contamination, (determined visually as instances where contigs had

526 anomalous GC content or tetranucleotide frequency), suspicious contigs were removed.

527 Redundancy was also used as a way to flag any genomes that needed further curation. After

528 inspection, curated contigs were exported using the program 'anvi-summarize' and reads were

529 re-mapped to the cleaned version of assemblies. The cleaned genomes were processed again for

530 visualization in anvi'o to ensure no erroneous contigs were included. Mapping quality was

531 inspected visually using the Integrative Genomics Viewer (IGV) (Robinson et al. 2011) and

532 Tablet (Milne et al. 2013) and manual curation was undertaken using mapped read data. Manual

533 inspection was used to determine if a circular genome could be considered closed and complete.

534 All contigs shorter than 1000 bp were removed from the genomes that were not closed after final

535 curation, and anvi'o was used to assess final genome completeness and redundancy (Eren et al.

536 2021).

537 **Phylogenomic analyses**

538    To generate a comprehensive phylogeny of the SAR11 clade, a suite of high-quality

539 genomes were curated. Even with an abundance of metagenomes, the high diversity among

540 SAR11 populations makes constructing reliable MAGs currently unfeasible, so to ensure the

541 phylogeny was as robust as possible, only isolate genomes and SAGs were included. The final

542 set of 493 SAR11 genomes for phylogenetic reconstruction included 81 genomes sequenced in

543 this study, 25 previously published reference genomes, and 387 previously published single

30

544 amplified genomes (SAGs), in addition to 20 isolate genomes from the family Rhodobacteraceae

545 that were used as an outgroup (Supplemental Table 3). The majority of SAGs included were

546 equal to or greater than 85% complete according to checkM (Parks et al. 2015). However,

547 genomes of lower quality from subclades of SAR11 with no high-quality representatives were

548 included to produce a comprehensive phylogeny, for example SAR11 Ic genomes that ranged

549 from 56.0 to 93.7 percent completion were also included (Thrash et al. 2014) (Supplemental

550 Table 3). Both previously identified subgroup V and IV genomes were excluded from these

551 analyses as subgroup V is not considered to be within SAR11 and the inclusion of subgroup IV

552 has not been rigorously investigated and thus its relationship to the *Pelagibacterales* is

553 questionable (Thrash et al. 2011; Viklund et al. 2013; Haro-Moreno et al. 2020; Muñoz-Gómez

554 et al. 2022).

555       We compared two gene sets to determine the most appropriate genes to use for

556 phylogenetic reconstructions of the SAR11 clade. This included the bac120 gene set utilized by

557 GTDB-Tk to determine the bacteria guide tree, and a curated gene set of marker genes derived

558 from the 200-genes previously demonstrated to be best fit for the *Alphaproteobacteria*

559 (Muñoz-Gómez et al. 2019) (Supplemental Table 4). To curate the second gene set, we generated

560 a custom HMM profile for the 200 *Alphaproteobacteria* genes with a noise cutoff term of

561 $1\times10^{-20}$, ran the HMM profile on all genomes using the anvi'o program `anvi-run-hmms`, and

562 generated a presence-absence matrix of genes in this model across genomes using the program

563 `anvi-script-gen-hmm-hits-matrix-across-genomes`. After evaluating the model hits across the

564 genomes matrix, we removed the genes that occurred in less than 90% of the genomes or those

565 that were redundant in more than 2% of the genomes from the *Alphaproteobacteria* 200-gene

566 collection, which resulted in a new collection with 165 genes, which is referred to as the

567 `SAR11_165` throughout our study (Supplemental Table 4). To generate a concatenated

568 alignment of the genes of interest for downstream phylogenomic analyses, a custom HMM

569 source was generated that encompassed the SAR11_165 genes. The program

570 `anvi-get-sequences-for-hmm-hits ` with the custom HMM source was then implemented to

571 extract and align genes of interest. The program trimAL 1.3 (Capella-Gutiérrez et al. 2009) was

572 then used to remove all positions that were missing in more than 50% of the genomes.

573 Phylogenies were generated with IQ-Tree v2.1.2 (Minh et al. 2020) with the best fit model

574 (LG+F+R10) chosen using ModelFinder (Kalyaanamoorthy et al. 2017) and 1,000 ultrafast

575 bootstraps. Phylogenies were rerooted appropriately in FigTree, and exported in NEXUS format

576 with the options selected to "Save as currently displayed" and "Include Annotations (NEXUS &

577 JSON only)". Once exported, phylogenies were then compared using the package phytools

578 (Revell 2024) in R (R Development Core Team 2011).

579         Once the extended phylogeny was established, a subset of SAR11 genomes was used to

580 generate a pruned phylogeny with the SAR_165 gene set. For this, we first used PyANI

581 (Pritchard et al. 2016) to dereplicate all genomes using 95% gANI as a cutoff, then excluded

582 SAGs that did not share at least 90% gANI with a neighboring genome, and finally included 10

583 genomes from the GTDB that spanned 10 families from the order *Rhodospirillales* as an

584 outgroup prior to recomputing the final phylogenomic tree as described above. The 95% ANI

585 dereplication cutoff was chosen to avoid read splitting during competitive read recruitment and

586 for any clusters in which isolate genomes were available, they were chosen as preferred

587 representatives.

32

588

**Classification and nomenclature**

590     The extended phylogeny was used to define cohesive genetic clusters at the distal end of

591 the SAR11 tree. Single genomes that did not share at least 90% ANI with a neighboring genome

592 were not classified into genera.

593     To determine how taxonomic levels across the SAR11 lineage would compare using

594 relative evolutionary distance, we implemented this approach as previously described (Ramfelt et

595 al. 2024). Briefly, a domain-level phylogeny was first constructed using the GTDB-Tk

596 de_novo_workflow (Chaumeil et al. 2019) with SAR11 isolate and SAGs as well as

597 "p__Chloroflexota" as the outgroup. Marker genes were identified from the input genomes using

598 GTDB-Tk `identify`, and then aligned with GTDB-Tk `align` (using the "–skip_gtdb_refs"

599 flag). Finally, a tree was constructed using FastTree v2.1.10 (model WAG+GAMMA) (Price et

600 al. 2010), rooted with the Chloroflexota outgroup. This phylogeny was used as the input for the

601 'scale_tree' program in PhyloRank v0.1.11 (https://github.com/dparks1134/PhyloRank) to

602 convert branch lengths into relative evolutionary distance (RED). RED values of 0.77 and 0.92

603 were used to assess how they would align with family and genus-level lineages, respectively.

604 These values were based on the distribution of internal nodes within the SAR11 clade and values

605 used previously for other family and genus-level lineages (Parks et al. 2018).

606

**Read recruitment**

608     To assess the distribution of the newly described strains described in this study and put

609 them into context with previously sequenced genomes, we used a read recruitment approach with

33

610 globally distributed metagenomes. The SAR11 genomes included in this study were grouped into

611 clusters that shared 95% average nucleotide identity (ANI) or greater and representatives from

612 these 95% gANI groups were then used for read recruitment (n = 314, Supplemental Table 12).

613 Results from read recruitment were extrapolated for the other genomes included in each 95%

614 gANI group.

615      Metagenomes used for recruitment included those sequenced in Kāneʻohe Bay (Tucker et

616 al. 2024b), the environment from which the genomes were isolated. Only samples from sites

617 previously categorized as "nearshore" and "offshore" (Tucker et al. 2024b) were used here.

618 Additionally, globally distributed previously published metagenomes were also used including

619 those from TARA Oceans expeditions (Sunagawa et al. 2015), station ALOHA (Mende et al.

620 2017), GEOTRACERS cruises (Biller et al. 2018), the eastern coast of Japan (Kudo et al. 2018;

621 Yoshitake et al. 2021), Monterey Bay (Mueller et al. 2015), and the ocean sampling day program

622 (Kopf et al. 2015) (Supplemental Table 6 for a list of appropriate references and details regarding

623 metagenomes included).

624      Once metagenomes were chosen, raw reads were downloaded using 'prefetch' and

625 'fasterq-dump' in the SRA toolkit. We automated the quality filtering of metagenomes,

626 metagenomic read recruitment, and profiling of recruited reads using the program

627 anvi-run-workflow (Shaiber et al. 2020) with the `--workflow metagenomics` flag, which

628 implements snakemake (Köster and Rahmann 2012) recipes for standard analyses in anvi'o.

629 Briefly, this workflow identified and discarded the noisy sequencing reads in metagenomes using

630 the program `iu-filter-quality-minoche` (Eren et al. 2013b), used SAR11 genomes to

631 competitively recruit short reads from metagenomes using Bowtie2 (Langmead and Salzberg

34

632 2012) SAMtools (Li et al. 2009) using the program `anvi-profile`, and finally merge individual

633 profiles into an anvi'o merged profile database using the program `anvi-merge`. The resulting

634 anvi'o merged profile database included essential data, including genome coverages and

635 detection statistics across metagenomes, for our downstream analyses. For coverage, we

636 primarily used the 'mean coverage Q2Q3' statistic, which represents the interquartile average of

637 coverage values where, for any given genome, the lowest 25% and the highest 25% of individual

638 coverage values are trimmed prior to calculating the average coverage from the remaining data

639 points, and thus minimizing the impact of biases due to highly conserved or highly variable

640 regions in the final coverage estimates. Visualization of read recruitment data mapped according

641 to the phylogeny constructed was completed using the program `anvi-interactive` with the

642 `--manual` flag.

643

644 **Metagenome profile clustering**

645       We performed a cluster analysis of metagenomes based on genome detection values from

646 the read recruitment step using the k-means algorithm, where we determined the `k` by

647 identifying the elbow of the curve of within-cluster sum of square values for increasing values of

648 `k` using the R code shared by Delmont et al. (2019) at https://merenlab.org/data/sar11-saavs/.

649 The results of the clustering analysis were visualized using anvi'o. To investigate how similar

650 detection patterns of genomes within genome clusters were, in addition to how similar or distinct

651 patterns were between genome-clusters, we performed a non-metric multidimensional scaling

652 (NMDS) analysis using the vegan package in R. Any metagenomes with zero detection across all

35

653 genomes were removed. The NMDS results were visualized using ggplot2 and plotly and an

654 interactive plot was generated with ggplotly.

655

656 **Investigation of C-P lyase pathway**

657     All genomes included in the extended phylogeny (n=X) were searched using

658 `anvi-search-functions` for the key enzyme in the C-P lyase pathway (*phnJ*) to determine the

659 capacity among high-quality SAR11 genomes to utilize the pathway. The genes upstream and

660 downstream of this essential gene were extracted from all 57 genomes and a pangenome was

661 used to compare the presence and absence of other key genes in the pathway as well as the

662 synteny of this region of the genome. The *phnJ* phylogeny (Supplemental Fig. 6) does not reflect

663 the relationships among genomes as demonstrated by the SAR_165 phylogeny (Fig. 1), which is

664 further evidence that this gene is located on a genomic island as previously described

665 (Molina-Pardines et al. 2023).


# Acknowledgments
<small>666</small>

36

# Competing interests

The authors declare no competing interests.

# Data availability

The assembled sequence data for genomes reported here are available at FigShare at

https://doi.org/10.6084/m9.figshare.28087454.v1.

# References

Acker, M., S. L. Hogle, P. M. Berube, T. Hackl, A. Coe, R. Stepanauskas, S. W. Chisholm, and D. J. Repeta. 2022. Phosphonate production by marine microbes: Exploring new sources and potential function. Proc. Natl. Acad. Sci. U. S. A. **119**: e2113386119.

Bankevich, A. and others. 2012. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. J. Comput. Biol. **19**: 455–477.

Becker, J. W., S. L. Hogle, K. Rosendo, and S. W. Chisholm. 2019. Co-culture and biogeography of Prochlorococcus and SAR11. ISME J. **13**: 1506–1519.

Biller, S. J. and others. 2018. Marine microbial metagenomes sampled across space and time. Sci Data **5**: 180176.

Braakman, R. and others. 2024. Global niche partitioning of purine and pyrimidine cross-feeding among ocean microbes. bioRxiv. doi:10.1101/2024.02.09.579562

Brandon, M. L. 2006. High-throughput isolation of pelagic marine macteria from the coastal subtropical

692    Pacific Ocean. University of Hawaiʻi at Mānoa.

693  Brown, M. V. and others. 2012. Global biogeography of SAR11 marine bacteria. Mol. Syst. Biol. **8**: 595.

694  Callahan, B. J., P. J. McMurdie, M. J. Rosen, A. W. Han, A. J. A. Johnson, and S. P. Holmes. 2016.

695    DADA2: High-resolution sample inference from Illumina amplicon data. Nat. Methods **13**: 581–583.

696  Capella-Gutiérrez, S., J. M. Silla-Martínez, and T. Gabaldón. 2009. trimAl: a tool for automated

697    alignment trimming in large-scale phylogenetic analyses. Bioinformatics **25**: 1972–1973.

698  Carini, P., L. Steindler, S. Beszteri, and S. J. Giovannoni. 2013. Nutrient requirements for growth of the

699    extreme oligotroph "Candidatus Pelagibacter ubique" HTCC1062 on a defined medium. ISME J. **7**:

700    592–602.

701  Carlson, C. A., R. Morris, R. Parsons, A. H. Treusch, S. J. Giovannoni, and K. Vergin. 2009. Seasonal

702    dynamics of SAR11 populations in the euphotic and mesopelagic zones of the northwestern Sargasso

703    Sea. ISME J. **3**: 283–295.

704  Chang, T., G. S. Gavelis, J. M. Brown, and R. Stepanauskas. 2024. Genomic representativeness and

705    chimerism in large collections of SAGs and MAGs of marine prokaryoplankton. Microbiome **12**:

706    126.

707  Chaumeil, P.-A., A. J. Mussig, P. Hugenholtz, and D. H. Parks. 2019. GTDB-Tk: a toolkit to classify

708    genomes with the Genome Taxonomy Database. Bioinformatics **36**: 1925–1927.

709  Craig, L., and J. Li. 2008. Type IV pili: paradoxes in form and function. Curr. Opin. Struct. Biol. **18**:

710    267–277.

711  Davies, T. J. 2006. Evolutionary ecology: when relatives cannot live together. Curr. Biol. **16**: R645–7.

712  Delmont, T. O. and others. 2018. Nitrogen-fixing populations of Planctomycetes and Proteobacteria are

713    abundant in surface ocean metagenomes. Nat Microbiol **3**: 804–813.

714  Delmont, T. O., E. Kiefl, O. Kilinc, O. C. Esen, I. Uysal, M. S. Rappé, S. Giovannoni, and A. M. Eren.

715    2019. Single-amino acid variants reveal evolutionary processes that shape the biogeography of a

716    global SAR11 subclade. Elife **8**. doi:10.7554/eLife.46497

717  Eiler, A., D. H. Hayakawa, M. J. Church, D. M. Karl, and M. S. Rappé. 2009. Dynamics of the SAR11

718    bacterioplankton lineage in relation to environmental conditions in the oligotrophic North Pacific

719    subtropical gyre. Environ. Microbiol. **11**: 2291–2300.

720  Eren, A. M. and others. 2021. Community-led, integrated, reproducible multi-omics with anvi'o. Nat

721    Microbiol **6**: 3–6.

722  Eren, A. M., Ö. C. Esen, C. Quince, J. H. Vineis, H. G. Morrison, M. L. Sogin, and T. O. Delmont. 2015.

723    Anvi'o: an advanced analysis and visualization platform for 'omics data. PeerJ **3**: e1319.

724  Eren, A. M., L. Maignien, W. J. Sul, L. G. Murphy, S. L. Grim, H. G. Morrison, and M. L. Sogin. 2013a.

725    Oligotyping: Differentiating between closely related microbial taxa using 16S rRNA gene data.

726    Methods Ecol. Evol. **4**: 1111–1119.

727  Eren, A. M., J. H. Vineis, H. G. Morrison, and M. L. Sogin. 2013b. A filtering method to generate high

728    quality short reads using illumina paired-end technology. PLoS One **8**: e66643.

729  Evans, J. T., and V. J. Denef. 2020. To dereplicate or not to dereplicate? mSphere **5**.

730    doi:10.1128/mSphere.00971-19

731  Giovannoni, S. J. and others. 2005. Genome streamlining in a cosmopolitan oceanic bacterium. Science

732    **309**: 1242–1245.

733  Giovannoni, S. J. 2017. SAR11 Bacteria: The Most Abundant Plankton in the Oceans. Ann. Rev. Mar.

734    Sci. **9**: 231–255.

735  Giovannoni, S. J., T. B. Britschgi, C. L. Moyer, and K. G. Field. 1990. Genetic diversity in Sargasso Sea

736    bacterioplankton. Nature **345**: 60–63.

737  Giovannoni, S. J., J. Cameron Thrash, and B. Temperton. 2014. Implications of streamlining theory for

738    microbial ecology. ISME J. **8**: 1553–1565.

739  Grote, J., J. C. Thrash, M. J. Huggett, Z. C. Landry, P. Carini, S. J. Giovannoni, and M. S. Rappé. 2012.

39

740     Streamlining and core genome conservation among highly divergent members of the SAR11 clade.

741     MBio **3**. doi:10.1128/mBio.00252-12

742 Haro-Moreno, J. M. and others. 2020. Ecogenomics of the SAR11 clade. Environ. Microbiol. **22**:

743     1748–1763.

744 Hug, L. A. and others. 2016. A new view of the tree of life. Nat Microbiol **1**: 16048.

745 Jain, C., L. M. Rodriguez-R, A. M. Phillippy, K. T. Konstantinidis, and S. Aluru. 2018. High throughput

746     ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. Nat. Commun. **9**: 5114.

747 Kalyaanamoorthy, S., B. Q. Minh, T. K. F. Wong, A. von Haeseler, and L. S. Jermiin. 2017. ModelFinder:

748     fast model selection for accurate phylogenetic estimates. Nat. Methods **14**: 587–589.

749 Kiefl, E., O. C. Esen, S. E. Miller, K. L. Kroll, A. D. Willis, M. S. Rappé, T. Pan, and A. M. Eren. 2023.

750     Structure-informed microbial population genetics elucidate selective pressures that shape protein

751     evolution. Sci Adv **9**: eabq4632.

752 Kopf, A. and others. 2015. The ocean sampling day consortium. Gigascience **4**: 27.

753 Köster, J., and S. Rahmann. 2012. Building and documenting workflows with python-based Snakemake.

754     GCB 49–56.

755 Kudo, T. and others. 2018. Seasonal changes in the abundance of bacterial genes related to

756     dimethylsulfoniopropionate catabolism in seawater from Ofunato Bay revealed by metagenomic

757     analysis. Gene **665**: 174–184.

758 Langmead, B., and S. L. Salzberg. 2012. Fast gapped-read alignment with Bowtie 2. Nat. Methods **9**:

759     357–359.

760 Li, H. and others. 2009. The Sequence Alignment/Map format and SAMtools. Bioinformatics **25**:

761     2078–2079.

762 López-Pérez, M., J. M. Haro-Moreno, F. H. Coutinho, M. Martinez-Garcia, and F. Rodriguez-Valera.

763     2020. The evolutionary success of the marine bacterium SAR11 analyzed through a metagenomic

40

764    perspective. mSystems **5**. doi:10.1128/mSystems.00605-20

765  Mende, D. R., J. A. Bryant, F. O. Aylward, J. M. Eppley, T. Nielsen, D. M. Karl, and E. F. DeLong. 2017.

766    Environmental drivers of a microbial genomic transition zone in the ocean's interior. Nat Microbiol

767    **2**: 1367–1373.

768  Milne, I., G. Stephen, M. Bayer, P. J. A. Cock, L. Pritchard, L. Cardle, P. D. Shaw, and D. Marshall. 2013.

769    Using Tablet for visual exploration of second-generation sequencing data. Brief. Bioinform. **14**:

770    193–202.

771  Minh, B. Q., H. A. Schmidt, O. Chernomor, D. Schrempf, M. D. Woodhams, A. von Haeseler, and R.

772    Lanfear. 2020. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the

773    Genomic Era. Mol. Biol. Evol. **37**: 1530–1534.

774  Molina-Pardines, C., J. M. Haro-Moreno, and M. López-Pérez. 2023. Phosphate-related genomic islands

775    as drivers of environmental adaptation in the streamlined marine alphaproteobacterial HIMB59.

776    mSystems **8**: e0089823.

777  Monaghan, E. A., K. C. Freel, and M. S. Rappé. 2020. Isolation of SAR11 Marine Bacteria from

778    Cryopreserved Seawater. mSystems **5**. doi:10.1128/mSystems.00954-20

779  Morris, R. M., M. S. Rappé, S. A. Connon, K. L. Vergin, W. A. Siebold, C. A. Carlson, and S. J.

780    Giovannoni. 2002. SAR11 clade dominates ocean surface bacterioplankton communities. Nature

781    **420**: 806–810.

782  Morris, R. M., K. L. Vergin, J.-C. Cho, M. S. Rappé, C. A. Carlson, and S. J. Giovannoni. 2005. Temporal

783    and spatial response of bacterioplankton lineages to annual convective overturn at the Bermuda

784    Atlantic Time-series Study site. Limnol. Oceanogr. **50**: 1687–1696.

785  Mueller, R. S. and others. 2015. Metagenome sequencing of a coastal marine microbial community from

786    monterey bay, california. Genome Announc. **3**. doi:10.1128/genomeA.00341-15

787  Muñoz-Gómez, S. A., S. Hess, G. Burger, B. F. Lang, E. Susko, C. H. Slamovits, and A. J. Roger. 2019.

41

788    An updated phylogeny of the Alphaproteobacteria reveals that the parasitic Rickettsiales and

789    Holosporales have independent origins. Elife **8**. doi:10.7554/eLife.42535

790 Muñoz-Gómez, S. A., E. Susko, K. Williamson, L. Eme, C. H. Slamovits, D. Moreira, P. López-García,

791    and A. J. Roger. 2022. Site-and-branch-heterogeneous analyses of an expanded dataset favour

792    mitochondria as sister to known Alphaproteobacteria. Nat. Ecol. Evol. **6**: 253–262.

793 Olm, M. R., A. Crits-Christoph, S. Diamond, A. Lavy, P. B. Matheus Carnevali, and J. F. Banfield. 2020.

794    Consistent metagenome-derived metrics verify and delineate bacterial species boundaries. mSystems

795    **5**. doi:10.1128/mSystems.00731-19

796 Pachiadaki, M. G. and others. 2019. Charting the complexity of the marine microbiome through

797    single-cell genomics. Cell **179**: 1623–1635.e11.

798 Paoli, L. and others. 2022. Biosynthetic potential of the global ocean microbiome. Nature **607**: 111–118.

799 Parks, D. H., M. Chuvochina, C. Rinke, A. J. Mussig, P.-A. Chaumeil, and P. Hugenholtz. 2022. GTDB:

800    an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank

801    normalized and complete genome-based taxonomy. Nucleic Acids Res. **50**: D785–D794.

802 Parks, D. H., M. Chuvochina, D. W. Waite, C. Rinke, A. Skarshewski, P.-A. Chaumeil, and P. Hugenholtz.

803    2018. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree

804    of life. Nat. Biotechnol. **36**: 996–1004.

805 Parks, D. H., M. Imelfort, C. T. Skennerton, P. Hugenholtz, and G. W. Tyson. 2015. CheckM: assessing

806    the quality of microbial genomes recovered from isolates, single cells, and metagenomes. Genome

807    Res. **25**: 1043–1055.

808 Price, M. N., P. S. Dehal, and A. P. Arkin. 2010. FastTree 2 – Approximately Maximum-Likelihood Trees

809    for Large Alignments. PLoS One **5**: e9490.

810 Pritchard, L., R. H. Glover, S. Humphris, J. G. Elphinstone, and I. K. Toth. 2016. Genomics and

811    taxonomy in diagnostics for food security: soft-rotting enterobacterial plant pathogens. Anal.

42

812      Methods **8**: 12–24.

813 Quast, C., E. Pruesse, P. Yilmaz, J. Gerken, T. Schweer, P. Yarza, J. Peplies, and F. O. Glöckner. 2013.

814      The SILVA ribosomal RNA gene database project: improved data processing and web-based tools.

815      Nucleic Acids Res. **41**: D590–6.

816 Ramfelt, O., K. C. Freel, S. J. Tucker, O. D. Nigro, and M. S. Rappé. 2024. Isolate-anchored comparisons

817      reveal evolutionary and functional differentiation across SAR86 marine bacteria. ISME J. **18**.

818      doi:10.1093/ismejo/wrae227

819 Rappé, M. S., S. A. Connon, K. L. Vergin, and S. J. Giovannoni. 2002. Cultivation of the ubiquitous

820      SAR11 marine bacterioplankton clade. Nature **418**: 630–633.

821 R Development Core Team, R. 2011. R: A Language and Environment for Statistical Computing,.

822 Revell, L. J. 2024. phytools 2.0: an updated R ecosystem for phylogenetic comparative methods (and

823      other things). PeerJ **12**: e16505.

824 Robinson, J. T., H. Thorvaldsdóttir, W. Winckler, M. Guttman, E. S. Lander, G. Getz, and J. P. Mesirov.

825      2011. Integrative genomics viewer. Nat. Biotechnol. **29**: 24–26.

826 Sanford, R. A., K. G. Lloyd, K. T. Konstantinidis, and F. E. Löffler. 2021. Microbial taxonomy run amok.

827      Trends Microbiol. **29**: 394–404.

828 Schattenhofer, M., B. M. Fuchs, R. Amann, M. V. Zubkov, G. A. Tarran, and J. Pernthaler. 2009.

829      Latitudinal distribution of prokaryotic picoplankton populations in the Atlantic Ocean. Environ.

830      Microbiol. **11**: 2078–2093.

831 Schwalbach, M. S., H. J. Tripp, L. Steindler, D. P. Smith, and S. J. Giovannoni. 2010. The presence of the

832      glycolysis operon in SAR11 genomes is positively correlated with ocean productivity. Environ.

833      Microbiol. **12**: 490–500.

834 Shaiber, A. and others. 2020. Functional and genetic markers of niche partitioning among enigmatic

835      members of the human oral microbiome. Genome Biol. **21**: 292.

43

836  Sosa, O. A., D. J. Repeta, E. F. DeLong, M. D. Ashkezari, and D. M. Karl. 2019. Phosphate-limited ocean

837      regions select for bacterial populations enriched in the carbon-phosphorus lyase pathway for

838      phosphonate degradation. Environ. Microbiol. **21**: 2402–2414.

839  Sunagawa, S. and others. 2015. Ocean plankton. Structure and function of the global ocean microbiome.

840      Science **348**: 1261359.

841  Sun, J. and others. 2016. The abundant marine bacterium Pelagibacter simultaneously catabolizes

842      dimethylsulfoniopropionate to the gases dimethyl sulfide and methanethiol. Nat Microbiol **1**: 16065.

843  Sun, J., L. Steindler, J. C. Thrash, K. H. Halsey, D. P. Smith, A. E. Carter, Z. C. Landry, and S. J.

844      Giovannoni. 2011. One carbon metabolism in SAR11 pelagic marine bacteria. PLoS One **6**: e23973.

845  Suzuki, M. T., O. Béjà, L. T. Taylor, and E. F. Delong. 2001. Phylogenetic analysis of ribosomal RNA

846      operons from uncultivated coastal marine bacterioplankton. Environ. Microbiol. **3**: 323–331.

847  Thrash, J. C. and others. 2011. Phylogenomic evidence for a common ancestor of mitochondria and the

848      SAR11 clade. Sci. Rep. **1**: 13.

849  Thrash, J. C., B. Temperton, B. K. Swan, Z. C. Landry, T. Woyke, E. F. DeLong, R. Stepanauskas, and S.

850      J. Giovannoni. 2014. Single-cell enabled comparative genomics of a deep ocean SAR11 bathytype.

851      ISME J. **8**: 1440–1451.

852  Tripp, H. J., J. B. Kitner, M. S. Schwalbach, J. W. H. Dacey, L. J. Wilhelm, and S. J. Giovannoni. 2008.

853      SAR11 marine bacteria require exogenous reduced sulphur for growth. Nature **452**: 741–744.

854  Tsementzi, D. and others. 2016. SAR11 bacteria linked to ocean anoxia and nitrogen loss. Nature **536**:

855      179–183.

856  Tucker, S. J., K. C. Freel, A. M. Eren, and M. S. Rappe. 2024a. Habitat-specificity in SAR11 is associated

857      with a handful of genes under high selection. bioRxiv. doi:10.1101/2024.12.23.630198

858  Tucker, S. J., K. C. Freel, E. A. Monaghan, C. E. S. Sullivan, O. Ramfelt, Y. M. Rii, and M. S. Rappé.

859      2021. Spatial and temporal dynamics of SAR11 marine bacteria across a nearshore to offshore

44

860    transect in the tropical Pacific Ocean. PeerJ **9**: e12274.

861  Tucker, S. J., Y. M. Rii, K. C. Freel, K. Kotubetey, A. H. Kawelo, K. B. Winter, and M. S. Rappe. 2024b.

862    Sharp transitions in phytoplankton communities across estuarine to open ocean waters of the tropical

863    Pacific. bioRxiv. doi: 10.1101/2024.05.23.595464v1.

864  Tully, B. J., E. D. Graham, and J. F. Heidelberg. 2018. The reconstruction of 2,631 draft

865    metagenome-assembled genomes from the global oceans. Sci Data **5**: 170203.

866  Vergin, K. L. and others. 2013. High-resolution SAR11 ecotype dynamics at the Bermuda Atlantic

867    Time-series Study site by phylogenetic placement of pyrosequences. ISME J. **7**: 1322–1332.

868  Vergin, K. L., H. J. Tripp, L. J. Wilhelm, D. R. Denver, M. S. Rappé, and S. J. Giovannoni. 2007. High

869    intraspecific recombination rate in a native population of Candidatus pelagibacter ubique (SAR11).

870    Environ. Microbiol. **9**: 2430–2440.

871  Viklund, J., T. J. G. Ettema, and S. G. E. Andersson. 2012. Independent genome reduction and

872    phylogenetic reclassification of the oceanic SAR11 clade. Mol. Biol. Evol. **29**: 599–615.

873  Viklund, J., J. Martijn, T. J. G. Ettema, and S. G. E. Andersson. 2013. Comparative and phylogenomic

874    evidence that the alphaproteobacterium HIMB59 is not a member of the oceanic SAR11 clade. PLoS

875    One **8**: e78858.

876  Villarreal-Chiu, J. F., J. P. Quinn, and J. W. McGrath. 2012. The genes and enzymes of phosphonate

877    metabolism by bacteria, and their distribution in the marine environment. Front. Microbiol. **3**: 19.

878  Waite, D. W. and others. 2020. Proposal to reclassify the proteobacterial classes Deltaproteobacteria and

879    Oligoflexia, and the phylum Thermodesulfobacteria into four phyla reflecting major functional

880    capabilities. Int. J. Syst. Evol. Microbiol. **70**: 5972–6016.

881  Wang, Z., and M. Wu. 2013. A phylum-level bacterial phylogenetic marker database. Mol. Biol. Evol. **30**:

882    1258–1262.

883  Wick, R. R., L. M. Judd, C. L. Gorrie, and K. E. Holt. 2017. Unicycler: Resolving bacterial genome

45

884    assemblies from short and long sequencing reads. PLoS Comput. Biol. **13**: e1005595.

885  Wilhelm, L. J., H. J. Tripp, S. A. Givan, D. P. Smith, and S. J. Giovannoni. 2007. Natural variation in

886    SAR11 marine bacterioplankton genomes inferred from metagenomic data. Biol. Direct **2**: 27.

887  Yoshitake, K. and others. 2021. Development of a time-series shotgun metagenomics database for

888    monitoring microbial communities at the Pacific coast of Japan. Sci. Rep. **11**: 12222.

889  Zhao, J. and others. 2024. Promiscuous and unbiased recombination underlies the sequence-discrete

890    species of the SAR11 lineage in the deep ocean. bioRxiv. doi:10.1101/2024.10.30.621061

891

892

893

894

895

896

897

898

899

900

901

902

# Supplemental Figures

All supplemental figures are available on FigShare at:

https://doi.org/10.6084/m9.figshare.28087760.

**Supplemental Figure 1**. **Sampling sites used for high-throughput culturing experiments. (a)** Location of Oʻahu in the Hawaiian archipelago in relations to Station ALOHA approximately 100 km north. **(b)** Map of the embayment on the windward side of Oʻahu with sites included in the Kāneʻohe Bay Time-series with sites classified as 'nearshore' (orange text), 'transition' (black text), or 'offshore' (turquoise text). Site SR4 in gray from which seawater media was collected for the cultivation experiments is also indicated. Bathymetry lines are approximate. **(c)** Flowchart outlining the high throughput cultivation (HTC) experiments conducted in 2017 and 2018 leading to the isolation of hundreds of SAR11 cultures and 79 new SAR11 isolate genomes. **(d)** Schematic SAR11 phylogeny to indicate which samples harbored genomes from which subgroups.

48
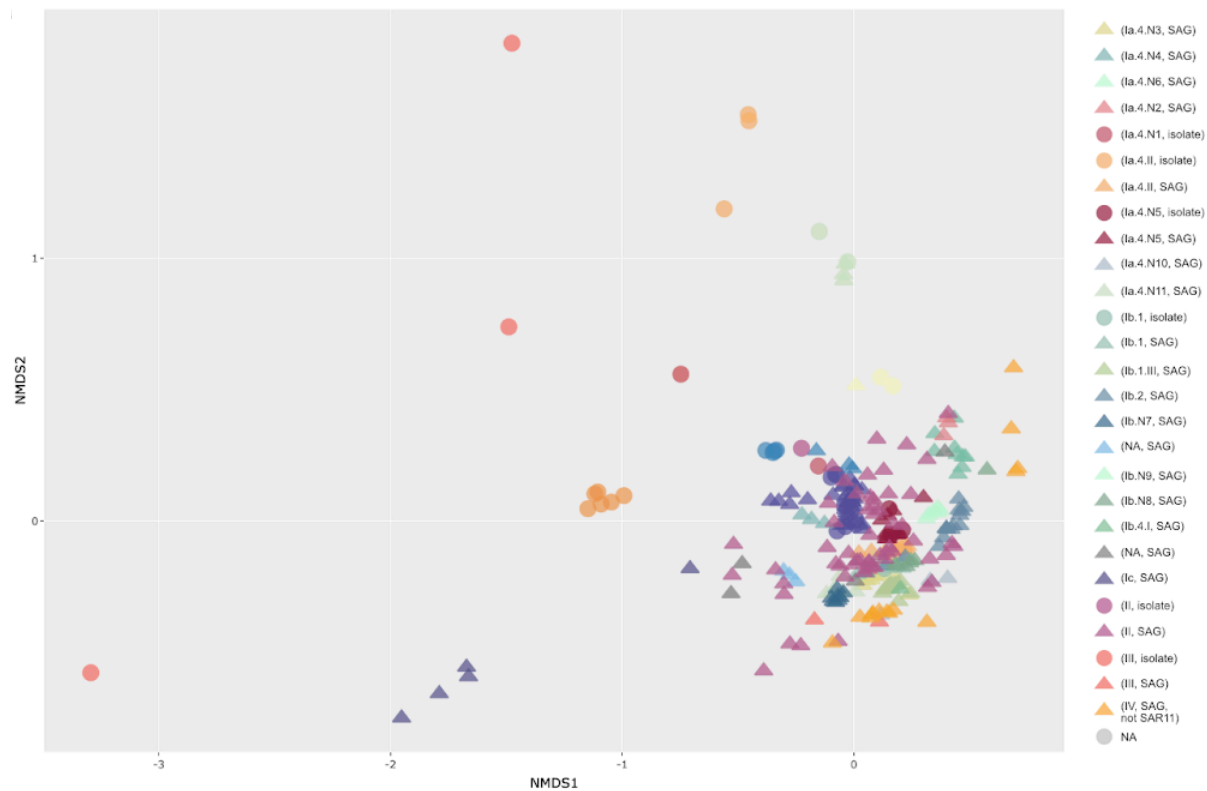
**Supplemental Figure 2. Phylogenomic tree of all 481 *Pelagibacterales* genomes initially included.** Of the 481 SAR11 genomes 106 were isolates and 375 SAGs and the phylogeny is based on a curated SAR11-specific set of 165 genes. Isolate origin is indicated and indicates if the genome was from this study, a previous isolate from Kāneʻohe Bay, or from another source. Genomes included in the pruned tree are also indicated.
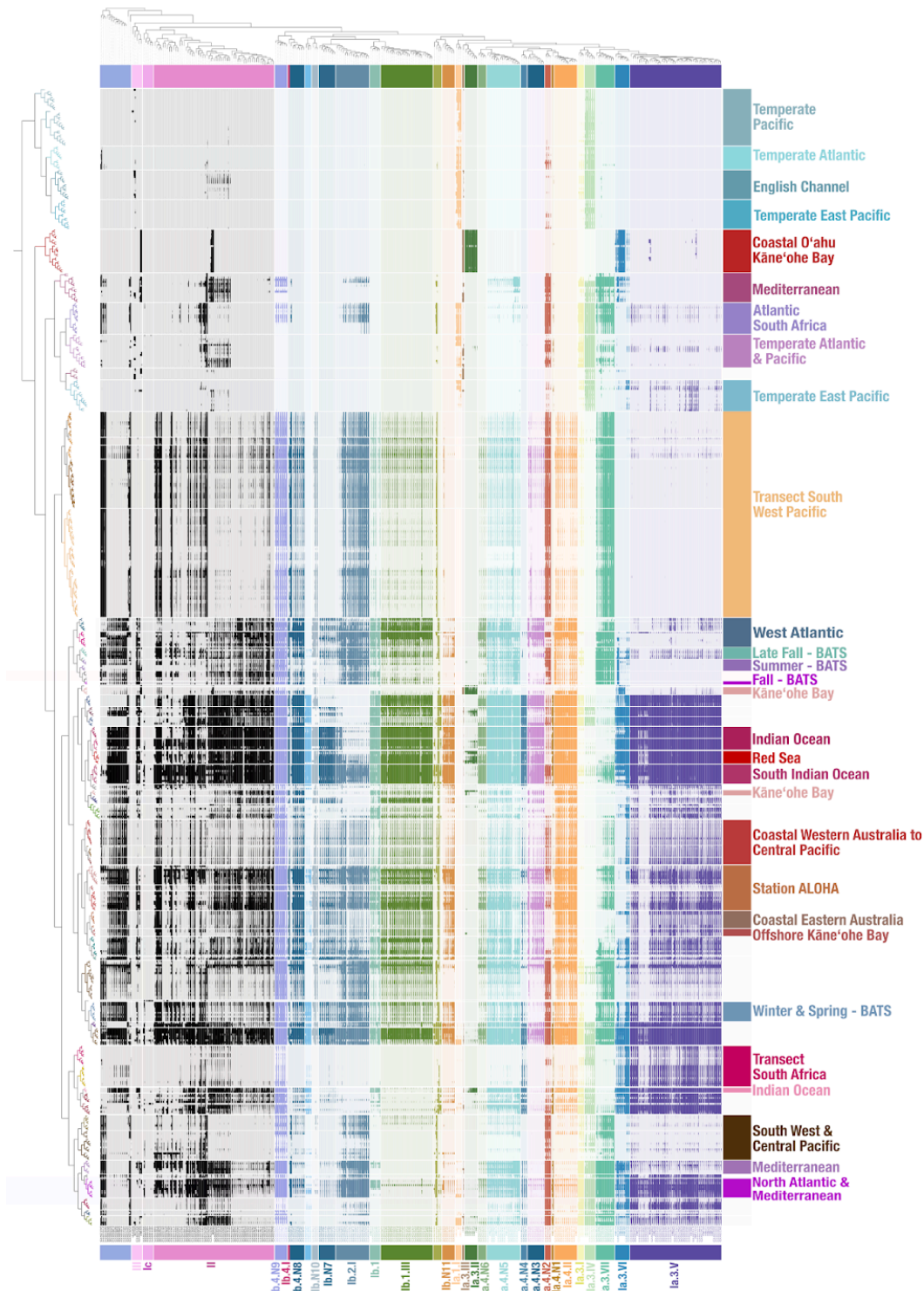
49

940

**Supplemental Figure 3. Pruned phylogenomic tree with 268 *Pelagibacterales* genomes.** Of the 268 genomes 50 were isolates and 218 SAGs based on a curated SAR11-specific set of 165 genes. Isolate origin is highlighted on the tree and indicates if the genome was from this study, a previous isolate from Kāneʻohe Bay, or from another source. Number of additional genomes in the same 95% gANI cluster are indicated as well by intensity of the bar which range from 0 to 45.
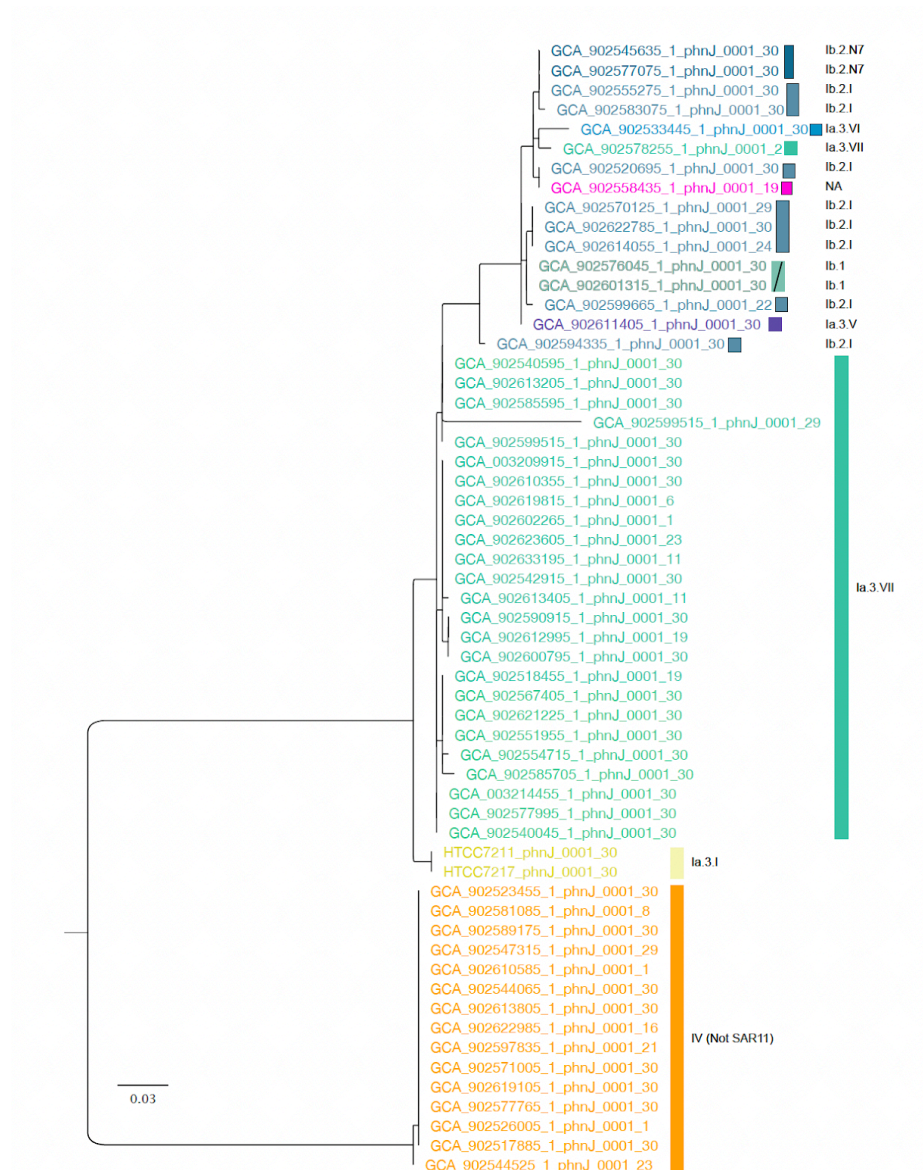
50

**Supplemental Figure 4. NMDS of genomes with detection data.** All of the genomes included in the analysis are included here, with distinction between isolate genome (circle) or SAG (triangle) indicated. Genomes not designated a subgroup noted as NA.

952

953

**Supplemental Figure 5. Read recruitment data including all metagenomes included in cluster analysis of metagenomes based on genome detection values using the k-means algorithm.** Subgroups are indicated at the bottom of the figure with labels indicating the metagenome groups along the right hand side of the figure.

52

958

**Supplemental Figure 6.** Phylogeny of the *phnJ* gene for all *Pelagibacterales* genomes in the pruned phylogeny data set.

960

961

# Supplemental Table Legends

963 All supplemental tables are available on FigShare at:

964 https://doi.org/10.6084/m9.figshare.28087490.v1.

965

**Supplemental Table 1. Summary of the 16S rRNA gene amplicon data from HTC17 and HTC18.**

968

**Supplemental Table 2. Detailed information for the isolate genomes reported in this study.** The genome summary information originated from checkM (v1.1.2). *Indicates the genome has been manually verified to be completely closed.

972

**Supplemental Table 3. Summary statistics for all genomes used for analyses in this study.** This includes isolates reported here, previously published isolate genomes, and high-quality single amplified genomes (SAGs) used in the extended SAR11 phylogeny. The genomes used in read-recruitment are indicated. *Indicates the accession is the IMG Genome ID not the NCBI Accession.

978

**Supplemental Table 4. Gene sets evaluated for use in SAR11 phylogenetics.** The sets evaluated include the bac120 (Parks et al., 2018) and a subset of 165 of the genes (SAR11_165) delineated for the Alphaproteobacteria (Wang and Wu 2013; Muñoz-Gómez, 2019).

982

**Supplemental Table 5. Summary of average genome statistics for the 23 genera established in the *Pelagibacteraceae* as well as the Ic, II, and III families.**

985

**Supplemental Table 6. Studies from which metagenomes were sourced.**

987

**Supplemental Table 7. List of all metagenomes used for read recruitment and accession numbers.**

990

**Supplemental Table 8. Detection values across genomes from all metagenomes used in read recruitment.**

993

**Supplemental  Table 9. Average detection across genome cluster for bins in Fig 3.**

995

**Supplemental Table 10. Coverage values across genomes from all metagenomes used in read recruitment.**

998

**Supplemental Table 11. Type genomes and classification hierarchy for the *Pelagibacterales* including proposed naming schemes.**

**Supplemental  Table 12. All of the final 95% ANI clusters defined including the cluster number, the final representative genome for that cluster and the list of other genomes in the same cluster.**