

1 **Novel pelagiphages prevail in the ocean**

2

3 Zefeng Zhang<sup>a, 1</sup>, Fang Qin<sup>a, 1</sup>, Feng Chen<sup>b</sup>, Xiao Chu<sup>c</sup>, Haiwei Luo<sup>c</sup>, Rui Zhang<sup>d</sup>, Sen

4 Du<sup>a</sup>, Zhen Tian<sup>a</sup>, Yanlin Zhao<sup>a, 2</sup>

5

6 <sup>a</sup>Fujian Provincial Key Laboratory of Agroecological Processing and Safety

7 Monitoring, College of Life Sciences, Fujian Agriculture and Forestry University,

8 Fuzhou, Fujian, China;

9 <sup>b</sup>Institute of Marine and Environmental Technology, University of Maryland Center 10

10 for Environmental Science, Baltimore, MD, USA;

11 <sup>c</sup>Simon F. S. Li Marine Science Laboratory, School of Life Sciences and State Key

12 Laboratory of Agrobiotechnology, The Chinese University of Hong Kong, Shatin,

13 Hong Kong SAR, China;

14 <sup>d</sup>State Key Laboratory of Marine Environmental Science, College of Ocean and Earth

15 Sciences, Institute of Marine Microbes and Ecospheres, Xiamen University, Xiamen,

16 Fujian, China;

17 <sup>1</sup>These authors contributed equally: Z.Z. and F.Q. contributed equally to this work.

18 <sup>2</sup>To whom correspondence may be addressed. Email: yanlinzhao@fafu.edu.cn.

19

20 **Keywords:** SAR11, pelagiphages, novel phage groups, Viromics.

## 21 **Abstract**

22 Viruses play a key role in biogeochemical cycling and host mortality, metabolism,  
23 physiology and evolution in the ocean. Viruses that infect the globally abundant  
24 marine SAR11 bacteria (pelagiphages) were reported to be an important component of  
25 the marine viral community. In this study, ten pelagiphages that infect three different  
26 *Pelagibacter* strains were isolated from various geographical locations and were  
27 genomically characterized. All ten pelagiphages are novel, representing four new  
28 lineages of the *Podoviridae* family. Although they share limited homology with  
29 cultured phage isolates, they are all closely related to some environmental viral  
30 fragments. Two HTVC023P-type pelagiphages are shown to be related to the  
31 abundant VC\_6 and VC\_8 viral populations of the Global Oceans Viromes (GOV)  
32 datasets. Interestingly, HTVC103P-type pelagiphages contain a structural module  
33 similar to that in SAR116 phage HMO-2011. Three HTVC111P-type pelagiphages  
34 and HTVC106P are also novel and related to GOV VC\_41 and VC\_67 viral  
35 populations, respectively. Remarkably, these pelagiphage represented phage groups  
36 are all globally distributed and predominant. Half of the top ten most abundant known  
37 marine phage groups are represented by pelagiphages. The HTVC023P-type group is  
38 the most abundant known viral group, exceeding the abundance of HTVC010P-type  
39 and HMO-2011-type groups. Furthermore, the HTVC023P-type group is also  
40 abundant throughout the water column. Altogether, this study has greatly broadened  
41 our understanding of pelagiphages regarding their genetic diversity, phage-host  
42 interactions and the distribution pattern. Availability of these newly isolated  
43 pelagiphages and their genome sequences will allow us to further explore their phage-  
44 host interactions and ecological strategies.

45

## 46 **Introduction**

47 As the most abundant biological entities in the ocean, viruses play critical roles  
48 in impacting marine biogeochemical cycling and shaping the microbial community  
49 structure and function (1-3). They also harbor enormous genetic diversity and diverse

50 metabolic potentials (4-7). Despite the fundamental importance of marine viruses, we  
51 just began to understand the diversity of marine viral communities. In the most recent  
52 decade, culture-independent metagenomic surveys (4-8), metagenomics fosmid (9,  
53 10) and single-cell genomics (SCGs) (11-13) have been used to explore the marine  
54 viral genetic and functional diversity and obtain novel viral genomic fragments from  
55 uncultivated viruses. For example, 15,280 viral populations belonging to 867 genus-  
56 level viral clusters were identified from the analysis of the Global Oceans Viromes  
57 (GOV) (5), 488,130 viral populations were recently identified from GOV 2.0 datasets  
58 (6), and more than 1,000 viral genomic fragments were retrieved from a single fosmid  
59 library from the Mediterranean deep chlorophyll maximum (MedDCM) (9). These  
60 studies unveiled enormous diversity of viruses in the ocean. In contrast, culture-  
61 dependent viral isolation has a limited contribution to reveal the marine viral diversity  
62 due to the fact that many bacterial groups are not easy to be cultivated in the  
63 laboratory. The number of cultured viruses from the ocean is far less compared to the  
64 number of omic-assembled viruses. The culture- independent studies have unveiled  
65 enormous diversity of viruses in the ocean, and at the same time, raised challenges on  
66 finding their potential hosts and unravelling their ecological and biological roles.  
67 Considerable efforts have been made to predict the hosts of viral sequences and to  
68 predict potential phage-host interactions (5, 9, 14). Despite these efforts, hosts of most  
69 viral clusters identified from the GOV still remain unknown and the majority of the  
70 viral clusters identified from the GOV lack any cultivated representative (5).

71         Although current isolated phages are insufficient for elucidating the natural viral  
72 diversity in the ocean, the isolation and genomic analysis of some important marine  
73 phages, such as cyanophages, SAR116 phage (15) and SAR11 phages (referred as  
74 pelagiphages) (16), have greatly facilitated the interpretation of marine virome  
75 datasets. The discovery of SAR11 phages (16) and a SAR116 phage (15) advocates  
76 for the importance of viral isolation. It was estimated that the isolation of  
77 pelagiphages and SAR116 phage HMO-2011 increased the number of known reads of  
78 in viral metagenomes by 30% (17). Therefore, isolation and sequencing more phages

79 infecting ecologically important bacterial hosts is of urgency in the area of marine  
80 viral ecology. In addition, having phage isolates in culture has the advantages of  
81 obtaining full genome sequences, gaining phage biological information, establishing  
82 cultivated virus-host model systems for a better understanding the phage infection  
83 process and ecological functions in marine ecosystems.

84 The order *Pelagibacterales* (SAR11) within the *Alphaproteobacteria*, is  
85 ubiquitous in the marine environments, accounting for approximately one-third of the  
86 oceanic prokaryotic cells (18-20), making the SAR11 clade the largest population of  
87 closely related heterotrophic bacteria on Earth. SAR11 bacteria are typical marine  
88 “difficult-to-culture” oligotrophic bacteria, exhibiting a slow growth rate and  
89 requiring unusual culturing condition (18, 21). Due to the difficulty in SAR11 bacteria  
90 culturing and pelagiphage isolation, the genomic and ecological study of pelagiphages  
91 has just begun to be addressed by a few studies. Pelagiphages are among the most  
92 abundant marine known phage groups, and they influence the population dynamics  
93 and evolution of SAR11 (16). Given the ecological importance of SAR11,  
94 pelagiphage has received much research attention since the four pelagiphage isolates  
95 were reported in 2013. Currently, 15 pelagiphages belonging to three distinct phage  
96 groups have been reported, and they possess diverse genetic contents and novel life  
97 strategies (16, 22). Efforts are still needed to further explore the diversity of  
98 pelagiphages in order to better understand their genomic evolution and phage-host  
99 interactions.

100 In this study, we isolated and sequenced ten new pelagiphages from diverse  
101 marine environments. High levels of genetic diversity were revealed across these  
102 pelagiphage genomes. We showed that these pelagiphages belong to four novel viral  
103 groups and are closely related to many environmental viral fragments. Finally,  
104 metagenomic recruitment analyses reveal the dominance of certain pelagiphage  
105 represented phage groups in the upper ocean as well as in the deep ocean.

106

107 **Results and discussion**

108 **Morphology and general features of newly isolated pelagiphages.** In this study, a  
109 culture-dependent approach was employed to further explore the genetic diversity of  
110 pelagiphages. Ten pelagiphages infecting three SAR11 strains (98.9-99.4% 16S rRNA  
111 gene sequence identity) that belong to the SAR11-Ia subclade were isolated from  
112 diverse marine environments (Table 1). All pelagiphages belong to the *Podoviridae*  
113 family with capsid length ranging from 55 to 69 nm (Fig. 1A). It is noteworthy that  
114 among 15 previously isolated pelagiphages (16, 22), 14 are podoviruses, and only one  
115 is myovirus. No siphovirus infecting SAR11 has been reported yet. At present, all  
116 reported pelagiphages infect closely related SAR11 Ia isolates. More novel phage  
117 groups are expected to be discovered when diverse SAR11 strains are used for phage  
118 isolation.

119 All newly isolated pelagiphage genomes were assembled into a circular contig,  
120 indicating their genome completeness. The general features of these pelagiphages are  
121 shown in Table 1. The genome sizes of these pelagiphages vary from 32.4 to 60.8 kb  
122 with a G+C content of 30.4% to 35.0%, which is close to the G+C content of their  
123 hosts (29.0 to 29.7%) and other previously reported pelagiphages (29.7 to 35.5%) (16,  
124 22). No tRNA sequences were identified in all 10 pelagiphage genomes.

125

126 **Marine pelagiphages possess genetic diversity and novelty.** Overall, the 10  
127 pelagiphages reported here share limited sequence homology with other cultured  
128 phage isolates. Comparative genomics analysis categorized these pelagiphages into  
129 four distinct phage groups (at the genus-level approximately) (Figs. 1B and 2). To  
130 date, seven distinct phage groups were identified from pelagiphage isolates, six of  
131 which belong to the *Podoviridae* family, suggesting that podoviruses exert primary  
132 top-down control on the abundance and dynamics of SAR11 population. The  
133 prevalence and dominance of Pelagibacter podoviruses in the ocean is supported by  
134 the dominance of their viromic matches (see later Viromic fragment recruitment  
135 analyses). Integrase genes and other genes related to the lysogenic life cycle were not  
136 identified in all 10 pelagiphage genomes, indicating that these pelagiphages all infect

137 SAR11 bacteria using the lytic infection strategy.

138

139 **Match to the most abundant GOV viral clusters.** HTVC023P and HTVC027P are  
140 closely related, belonging to a novel HTVC023P-type phage group (Fig. 2A). Thirty-  
141 seven ORFs are shared between HTVC023P and HTVC027P (with 31 to 93% amino  
142 acid identity) and they have a conserved overall genome arrangement with few gene  
143 rearrangements (Fig. 2A). HTVC023P and HTVC027P exhibit no significant genomic  
144 synteny with other known phage isolates, thus possessing genomic and evolutionary  
145 novelty. Approximately half of the predicted ORFs from both HTVC023P and  
146 HTVC027P genomes show homology to genes from other types of phage and  
147 bacterial genomes. The remaining ORFs have no homologs in the NCBI-RefSeq  
148 database and some only hit environmental sequences. Of the predicted ORFs, only  
149 17% were assigned to putative biological functions based on the sequence homology  
150 analysis.

151 ORFs encoding the proteins necessary for phage DNA replication, packaging,  
152 morphology and lysis were identified. Although HTVC023P and HTVC027P  
153 resemble podoviruses in morphology, only few ORFs in their genomes have the best  
154 hit to other phages in the *Podoviridae* family. In the DNA replication region, both  
155 HTVC023P-type pelagiphages contain a DNA polymerase gene, a DNA helicase gene  
156 and a few function-unknown ORFs that are homologous to genes from siphoviruses,  
157 including *Dinoroseobacter* phage vB\_DshS-R5C, *Proteobacterial* phage phiJL001  
158 and several *Yuavirus* siphophages, with low amino acid identity (ranging from 26 to  
159 39%). Very few putative structural proteins could be identified with weak homology  
160 to other phage structural proteins. For example, HTVC027P ORF85 is homologous to  
161 the putative structure protein from *Cellulophaga* phage 18:3, and two others ORFs  
162 (HTVC023P ORF69 and HTVC027P ORF69) have small regions of homology to the  
163 putative tail fiber genes from some phage genomes. These results imply the existence  
164 of a novel set of phage structural proteins in the HTVC023P and HTVC027P  
165 genomes. Terminase large and small subunit (TerL and TerS) genes involved in DNA

166 packaging were predicted, which are more closely related to TerL genes in some  
167 bacterial genomes. Both HTVC023P and HTVC027P harbor a GroES gene that  
168 encode a 10 kDa co-chaperonin. Many bacteria contain the GroEL/GroES molecular  
169 chaperonin system that is responsible for proper folding of many proteins, thus  
170 playing an important role in cell growth and cellular phage assembly (23, 24).  
171 Bacteriophage encoded cochaperonins were identified and studied in some phage  
172 genomes (25, 26). Both HTVC023P and HTVC027P GroES genes do not show  
173 significant homology to any phage GroES, while being mostly related to GroES  
174 sequences retrieved from marine viromes (27). Sequence analyses suggests that  
175 GroES in HTVC023P and HTVC027P are clustered with GroES clusters 14 and  
176 cluster 1, respectively, which were among the most abundant GroES clusters  
177 identified from viromic datasets (27).

178 The result of the DNA polymerase gene phylogeny analysis shows that  
179 HTVC023P and HTVC027P DNA polymerases are placed with clade III DNAP  
180 genes identified in a previous shotgun metaviromes study (28), and DNA polymerases  
181 from vB\_DshS-R5C, *Yuavirus* siphoviruses and phiJL001 are more distantly related  
182 (Fig. 1C). Clade III DNA polymerases accounted for 77% of all identified DNA  
183 polymerases from the Chesapeake Bay, Gulf of Maine and Dry Tortugas (28). It was  
184 previously speculated that Clade III DNA polymerases are likely from lysogenic  
185 phages (28), whereas our study shows that this DNA polymerase group is related to  
186 lytic podoviruses represented by HTVC023P-type pelagiphages.

187 Gene-content-based network analysis reveals that 443 viral sequences ( $\geq 20$ kb)  
188 from diverse ocean regions were grouped into a viral cluster (VC\_009) with  
189 HTVC023P-type pelagiphages (Fig. 3), suggesting that that close relatives of  
190 HTVC023P-type pelagiphages exhibit globally distribution pattern. Phylogenetic  
191 analysis based on the VC\_009 DNA polymerase sequences reveals a high level of  
192 diversity (Fig. 4). We notice that GOV populations grouped with HTVC023P-type  
193 pelagiphages are exclusively from GOV viral clusters VC\_6 and VC\_8 (5). VC\_6 and  
194 VC\_8 were two of the most globally abundant viral clusters identified in GOV study



195 (5). Genomic analysis reveals the genetic relatedness and genome synteny between  
196 HTVC023P-type pelagiphages and representative contigs from GOV VC\_6 and  
197 VC\_8, showing a high degree of synteny (Fig. 5A). Approximately half of the contigs  
198 in GOV VC\_6 and VC\_8 share more than 40% genes with HTVC023P-type  
199 pelagiphages and most of the remaining contigs share more than 20% genes with  
200 HTVC023P-type pelagiphages. In most cases, the low percentage of shared genes  
201 between viral populations and HTVC023P-type pelagiphages is due to some  
202 environmental viral sequences covering the nonconserved variable phage genome  
203 regions (data not shown). These results suggest that HTVC023P-type pelagiphages  
204 and most viral populations from GOV VC\_6 and VC\_8 can be grouped at the  
205 genus/subfamily-level. The phylogenetic analysis reveals that all GOV VC\_6 and  
206 VC\_8 DNA polymerases are clustered with HTVC023P-type DNA polymerases, with  
207 36% to 87% amino acid identity (*S1 Appendix*, Fig. S1).

208 These two pelagiphages also show high homology with a viral single-amplified  
209 genome contig, vSAG 37-F6 (13) (Over 80% of the predicted proteins in vSAG 37-  
210 F6) (Fig. 5A). The vSAG 37-F6 population was reported to be closely related to GOV  
211 VC\_6 and VC\_8 and has been shown to be abundant in several oceanic regions (13).  
212 SAR11 was recently predicted as putative host of vSAG 37-F6 population by single-  
213 cell genomics (29). However, before our study, phage-host system of these extremely  
214 important viral clusters still remained unavailable. It is noteworthy that in an earlier  
215 study, homologs of HTVC023P-type genomes were also found in single cell genomic  
216 analyses of *Verrucomicrobia* and *Bacteroidetes* (AAA160P02 and AAA164-I21) (*S1*  
217 *Appendix*, Fig. S2) (12), suggesting that members of the HTVC023P-type group may  
218 infect different taxonomic groups of bacteria. Further investigation based on culture-  
219 independent or culture-dependent studies are required to explore the diversity and  
220 infected hosts of this important viral group.

221 Taken together, close relatives of HTVC023P-type pelagiphages have been  
222 previously identified from some culture-independent studies and were revealed to be  
223 extremely abundant. These two HTVC023P-type pelagiphages are first known



224 cultured representatives of this important viral group.

225

226 **Homology to the HMO-2011-type phage group.** Three pelagiphages, HTVC103P,  
227 HTVC104P and HTVC115P are closely related, belonging to a novel HTVC103P-  
228 type group (Figs. 1B and 2B). Approximately 15% of their ORFs were assigned with  
229 putative functions. HTVC103P-type pelagiphages exhibit novel genomic  
230 architectures, containing two functional modules, including a DNA replication  
231 module and a phage structural and packaging module (Fig. 2B). In the DNA  
232 replication module, DNA polymerase, DNA primase and single-strand binding  
233 protein were predicted from all three HTVC103P-type pelagiphage genomes. The  
234 closest homologs of HTVC103P-type DNA polymerases and primases from isolated  
235 phages are those found in members of the *Autographivirinae* subfamily and  
236 Cobavirus group roseophages. Interestingly, HTVC103P and HTVC104P both share  
237 12 genes with SAR116 phage HMO-2011 and HTVC115P shares 16 genes with  
238 HMO-2011. Most of the HMO-2011 homologs are in the structural and packaging  
239 modules, including genes encoding capsid, portal and terminase (Fig. 2B). In addition,  
240 there is considerable amino acid identity and conserved gene synteny between  
241 HTVC103P-type pelagiphages and HMO-2011 in this region (27-68% amino acid  
242 identity). Gene content-based network analysis also reveals the relatedness between  
243 HTVC103P-type pelagiphages and HMO-2011-type phages (Fig. 3). These results  
244 suggest that HTVC103P-type genomes are probably composed of a DNA replication  
245 module and a phage structural and packaging module with distinct evolutionary  
246 origins and histories. These results indicate that horizontal gene exchange of the  
247 function module among phages may play an important role in driving evolution and  
248 genetic diversity of pelagiphages. It is likely that the transfer of a set of structural or  
249 DNA replication machinery genes occurred when divergent phages infected the same  
250 host cell or when there was contact between a resident prophage and an invading  
251 phage. Further investigation are required to illuminate the evolutionary trajectories of  
252 this novel phage group.

253 The DNA polymerase gene based phylogeny reveals that HTVC103P-type DNA  
254 polymerases are grouped with clade I DNA polymerases and are more distantly  
255 related to DNA polymerases from other *Autographivirinae* phages and Cobavirus  
256 roseophages (Fig. 1C). Clade I was another abundant DNA polymerase clade that was  
257 previously identified from shotgun metaviromes (28). In contrast, the structural genes  
258 of HTVC103P-type pelagiphages are most similar to those in HMO-2011-type  
259 phages. This result suggests that a phylogenetic approach based on a single gene has a  
260 limitation in revealing the evolutionary relationship among various phages. A gene  
261 content-based network can be used as a complement. Network analysis shows that a  
262 group of environmental viral fragments (49 sequences,  $\geq 20\text{kb}$ ) were clustered with  
263 HTVC103P-type pelagiphages, forming a viral cluster VC\_005 (Fig. 3). VC\_005  
264 shows distant relatedness to the HMO-2011-type group (Fig. 3). The DNA  
265 polymerase gene phylogeny reveals that the DNA polymerase genes of VC\_005 all  
266 cluster with HTVC103P-type pelagiphages and are distinct from HMO-2011-type  
267 DNA polymerases (*S1 Appendix*, Fig. S3).

268

269 **The new HTVC111P-type group and HTVC106P pelagiphage.** The HTVC111P-  
270 type phage group currently comprises pelagiphage HTVC111P, HTVC112P,  
271 HTVC026P and HTVC202P. Within this group, 45% to 67% of genes were shared  
272 (Fig. 1B). Genes responsible for phage replication, morphology, packaging, and lysis  
273 were identified from HTVC111P-type pelagiphage genomes with homology to genes  
274 from diverse bacteria and bacteriophages (Fig. 2C). The DNA polymerase gene was  
275 not found in the HTVC111P-type pelagiphage genomes and structural genes show  
276 very weak similarity to proteins from other known phage genomes, suggesting that  
277 this group of phages contain novel morphogenesis modules and may rely more on  
278 host replication system. The TerL gene in the HTVC111P-type genomes also show no  
279 significant similarity to other phage TerL genes.

280 Pelagiphage HTVC106P also lacks a clear relationship to any known phage  
281 isolates. Of 70 predicted ORFs in HTVC106P, approximately 40% have homologs in

282 other organisms in the NCBI-RefSeq database and only 11 ORFs were assigned with  
283 putative functions (Fig. 2D). For the remaining ORFs, most were highly similar to  
284 genes found only in metagenomic sequences. The ORFs assigned with functions are  
285 involved in DNA processing, virion morphogenesis, DNA packaging and lysis. The  
286 DNA replication genes were not identified in the HTVC106P genome, suggesting that  
287 HTVC106P may rely more on the host replication system or contain a novel  
288 replication module. The morphogenesis genes of HTVC106P show homology to some  
289 phages; for example, the HTVC106P portal protein, capsid and scaffolding protein  
290 exhibit distant homology with those of *Bruynoghevirus* phages, and the HTVC106P  
291 tail fiber protein and some other structural proteins are homologous with those in  
292 pelagiphage HTVC010P. HTVC106P contains a TerL related to the TerL in  
293 pelagiphage HTVC010P (49% amino acid identity), suggesting that HTVC106P  
294 possibly shares a conserved DNA packaging strategy with HTVC010P.

295 The close relatedness between HTVC111P-type pelagiphages, HTVC106P and  
296 some metagenomic viral fragments are also revealed by network analysis, with 165  
297 and 75 sequences ( $\geq 20$ kb) grouped with HTVC111P-type pelagiphages and  
298 HTVC106P, respectively (see VC\_016 and VC\_018 in Fig. 3). Phylogeny of capsid  
299 protein sequences reveals a high level of diversity of these two viral clusters (*SI*  
300 *Appendix*, Figs. S4 and S5). Network analysis also reveals the affiliation of  
301 pelagiphages with previously identified viral clusters in GOV. The GOV viral  
302 populations grouped with HTVC111P-type pelagiphages and HTVC106P are  
303 exclusively from GOV viral cluster VC\_41 and VC\_67, respectively (5).  
304 Approximately 80% of the viral populations ( $\geq 20$ kb) of VC\_41 and 70% of the viral  
305 populations ( $\geq 20$ kb) of VC\_67 are clustered with HTVC111P-type pelagiphages and  
306 HTVC106P, respectively. Genomic analysis reveals the genetic relatedness and  
307 genome synteny between HTVC111P-type pelagiphages and representative contigs  
308 from GOV VC\_41, HTVC106P and representative contigs from VC\_67 (Fig. 5B, 5C).  
309 The hosts of GOV VC\_41 and VC\_67 were not predicted yet (5). Most viral  
310 populations in GOV VC\_41 and VC\_67 share more than 20% genes with

311 HTVC111P-type phages and HTVC106P, respectively, indicating a possible  
312 relationship between these phages at the subfamily-level.

313

314 **Pelagiphage lytic and lysogenic developmental strategies.** Although the integrase  
315 gene was not identified from all ten pelagiphages (described above), a search for  
316 integrase genes reveals that a total of 27 environmental viral sequences closely related  
317 to HTVC106P were found to contain a tyrosine integrase gene (PFAM, PF00589) (*SI*  
318 *Appendix*, Fig. S6). Phage integrase mediates the site-specific recombination between  
319 phage sequence and bacterial sequences(30, 31). In addition, most of these integrase-  
320 containing viral sequences possess a sequence identical to SAR11 tRNA sequences  
321 (tRNA-Thr or tRNA-Asn) (*SI Appendix*, Fig. S6), which are likely to be the phage  
322 integration sites. In contrast, no identifiable integrases were found from  
323 environmental viral sequences closely related to other three phage groups. Although  
324 no intact prophages have been found in SAR11 genomes, lysogenic infection has been  
325 reported in HTVC019P-type pelagiphages (22). Our analysis suggests that a portion  
326 of the HTVC106P-type pelagiphages can infect the hosts via the lysogenic cycle,  
327 while other three pelagiphage types may have a strict lytic life strategy. These  
328 findings indicate that pelagiphage have diverse life strategies and lytic infection  
329 strategy is presumably the predominant form of pelagiphage-SAR11 interaction.

330

331 **Pelagiphages dominate the marine viromes.** The ecological importance of  
332 pelagiphages are reflected by their sheer abundance and ubiquity, as well as the  
333 predominance of their hosts in the ocean. In recent years, the marine viromic reads  
334 increased exponentially, providing valuable resources for assessing the relative  
335 abundance and distribution pattern of important viral groups in the ocean (5, 6, 32). In  
336 this study, a total of 174 marine virome datasets from various oceanic sites were  
337 downloaded for the viromic fragment recruitment analysis (*SI Appendix*, Table S1).  
338 We mainly used the reciprocal best-hit strategy to estimate the relative abundance of  
339 phage groups. Instead of accessing the relative abundance of each viral genome by

340 recruiting viromic reads with a high nucleotide identity (>95% or 90%) (9, 10, 13, 33,  
341 34), reciprocal recruitment method estimates the relative abundance of different phage  
342 groups at the genus or subfamily level. Due to the fact that a viral taxonomic group  
343 comprise evolutionarily diverse genotypes, reciprocal recruitment analysis can obtain  
344 an estimation of the relative abundance of important phage groups.

345 It is striking that within the top ten most abundant known phage groups at each  
346 marine viromic dataset, approximately half were represented by pelagiphages (Fig. 6  
347 and Dataset S1), suggesting that pelagiphages are abundant components of the marine  
348 viral communities. The HMO-2011-type group, some cyanophage- and roseophage-  
349 represented groups were also abundant throughout the world's ocean (Fig. 6). The  
350 new HTVC023P-type group was the most abundant viral group in the majority (80%)  
351 of the analyzed viromic datasets, exceeding the HMO-2011-type group and the  
352 HTVC010P-type group (Fig. 7). On average, the HTVC023P-type group is 2.8 and  
353 2.2 times more abundant than the HMO-2011-type group and the HTVC010P-type  
354 group, respectively. The reads assigned to the HTVC023P-type group accounted for  
355 0.55% to 7.59% of total viromic reads in various oceanic viromes, suggesting that the  
356 HTVC023P-type phage group dominates the marine viromes (Dataset S1).

357 The HTVC103P-type group ranked among the five most abundant viral groups  
358 in most viromic datasets. A remarkable feature of HTVC103P-type pelagiphages is  
359 that they harbor a set of structural genes that are homologous to those in HMO-2011-  
360 type phages. HMO-2011-type phages were highly represented in some ocean viromes,  
361 and the known host of HMO-2011-type phages currently comprise SAR116 and RCA  
362 roseobacters (15, 35). The structural gene homology between HTVC103P-type  
363 pelagiphages and HMO-2011-type phages raise the possibility that a portion of the  
364 viromic reads that were previously assigned to the HMO-2011-type group might be  
365 assigned to the HTVC103P-type group when HTVC103P-type genomes are included  
366 in the analysis. It was then estimated that the reads that were assigned to the HMO-  
367 2011-type group decreased by approximately 12% when HTVC103P-type  
368 pelagiphage genomes were included in the analysis (data not shown), thus

369 demonstrating that the HTVC103P-type group contributed to the abundance of HMO-  
370 2011-type phages in previous studies. When more related phage isolates are available,  
371 the analysis on phage relative abundance might be different to some extent.

372 Among all known pelagiphage groups, the HTVC112P-type and HTVC106P-  
373 type were less abundant than other pelagiphage groups, but they were still abundant  
374 and globally distributed. In the upper waters, the HTVC111P-type appeared to be as  
375 abundant as the HTVC019P-type group and the T4-like cyanomyoviruse group.

376 The above results further support that pelagiphages are extremely abundant and  
377 widely distributed in the ocean. Considering the ubiquity and dominance of SAR11  
378 bacteria in the ocean, they are prone to be attacked by viruses. The vast population  
379 sizes of SAR11 support the thriving of pelagiphages. The prevalence of pelagiphages  
380 in marine viromes suggests that SAR11 populations are under intense phage infection  
381 pressure. Podoviruses act as a primary contributor to SAR11 mortality and exert  
382 major control on the abundance and dynamics of SAR11 population. Although there  
383 is presumably significant cell loss of SAR11 populations due to the viral predation, as  
384 described by the King of the Mountain (KoM) hypothesis, the high recombination  
385 frequencies of SAR11 may also influence the distribution of phage-defense alleles,  
386 maintaining the coexistence of a high abundance of host and phages (16).

387 Among all known pelagiphage represented phage groups, the dominance of  
388 HTVC023P-type phages raises the question of what characteristics make them most  
389 successful. It presumably links to their biological traits; that is, they are likely to have  
390 higher infection efficiency (faster replicating) or they are more competitive when  
391 competing with other types of phages for hosts in complex viral assemblages.

392

393 **Vertical profiles of pelagiphages.** We compared the relative abundance of major  
394 marine phage groups at different water depths. There were significant variations in the  
395 relative abundance of most phage groups on the vertical scale (Fig. 6A, B). The  
396 relative abundance of the HTVC023P-type group appeared lower in mesopelagic (200  
397 to 1000 m) and bathypelagic (1000 to 4000 m) samples than in epipelagic samples.



398 Considering the cellular contamination reported in the bathypelagic viromes (5), it is  
399 possible that the HTVC023P-type exhibited comparable abundance throughout the  
400 water column. Remarkably, we observed that the HTVC023P-type group far exceeded  
401 other phage groups in the deep waters (200 to 4000 m) (Figs. 6 and 7). For example,  
402 in bathypelagic viromes, the HTVC023P-type group was on average 8 and 6 times  
403 more abundant than the HMO-2011-type and HTVC010P-type, respectively.  
404 Moreover, four of the five most abundant viral groups in mesopelagic and  
405 bathypelagic samples were represented by pelagiphages. In congruence, SAR11 is  
406 abundant throughout the oceanic water column and is abundant in coastal and open  
407 ocean (18, 36, 37); In deep waters, SAR11 subclades Ic, Iib, and Vb dominate the  
408 SAR11 populations (38, 39); thus, in the deep ocean, members of these phage groups  
409 are likely to infect these deep ocean SAR11 ecotypes. These results suggest that  
410 pelagiphages could be as important in deep ocean ecosystems as they are in the upper  
411 ocean. In contrast, the HMO-2011-type group only predominates the upper ocean,  
412 which is consistent with previous study (40). HMO-2011-type phages were found  
413 infecting SAR116 and RCA rosephages, which mainly occupy the upper ocean (15,  
414 35). Furthermore, we observed that the relative abundance of all pelagiphage groups  
415 did not exhibit significant variation between the coastal viromes and noncoastal  
416 viromes (*SI Appendix*, Fig. S7).

417

418 **Conclusions.** The 10 new pelagiphages that were described in this study greatly  
419 expand our current knowledge on the abundance, diversity, distribution and genomic  
420 evolution of viruses that infect SAR11 bacteria and further reinforce their ecological  
421 importance. Metagenomic recruitment analyses demonstrate that all these pelagiphage  
422 represented phage groups exhibit global distribution pattern and the HTVC023P-type  
423 group is the most dominant known viral group in the ocean. The predominance of  
424 these pelagiphages in marine viromes suggests that they could play an important role  
425 in controlling SAR11 population dynamics and influencing global carbon cycling.  
426 These new pelagiphages and their hosts will serve as useful model systems subjected



427 to the further investigations of various interactions between pelagiphages and their  
428 hosts, phage driven host evolution and dynamics, as well as the potential ecological  
429 impact of pelagiphages. It will also be interesting to study the mechanisms explaining  
430 the dominance of major phage groups. This study is another example of how phage  
431 isolation can improve the interpretation of marine viromic datasets, thus highlighting  
432 the irreplaceable power of culture-dependent phage isolation and cultivation in the  
433 study of marine virus functions and diversity. So far, all current known pelagiphage  
434 isolates were isolated from strains from the SAR11 Ia subclade. Future isolation of  
435 pelagiphages that infect other SAR11 subclades may reveal more novel phage lineages.

436

### 437 **Methods**

438 **SAR11 strains, media, and growth conditions.** SAR11 strains *Pelagibacter*  
439 HTCC7211 and *Pelagibacter* HTCC1062 were grown in an artificial seawater-based  
440 medium amended with 1 mM NH<sub>4</sub>Cl, 100 μM KH<sub>2</sub>PO<sub>4</sub>, 1 μM FeCl<sub>3</sub>, 100 μM  
441 pyruvate, 50 μM glycine, 50 μM methionine and excess vitamins (21). HTCC7211  
442 and HTCC1062 were grown at 17 °C and 20 °C, respectively. *Pelagibacter*  
443 FZCC0015 was isolated from Pingtan coast in 2017, detailed information on  
444 FZCC0015 has been described in earlier work (22). FZCC0015 was grown in a  
445 natural seawater-based medium amended with 100 μM pyruvate, 50 μM glycine, 50  
446 μM methionine and excess vitamins at 23 °C.

447

448 **Source waters and pelagiphage isolation.** Water samples were collected from a  
449 variety of oceanic sampling stations (Table 1). To obtain the cell-free fraction, the  
450 samples were filtrated through 0.1 μm-pore-size filters. The filtrates were stored in the  
451 dark at 4 °C until used for phage isolation. Isolation procedures for pelagiphages were  
452 described in detail previously (16, 22). Briefly, 0.1 μm filtered samples were  
453 inoculated with SAR11 cultures. Cell growth was monitored using a Guava EasyCyte  
454 cell counter (Millipore, Guava Technologies). When a decrease in cell densities was  
455 detected, the presence of phage particles was confirmed by epifluorescence

456 microscopy. Purified pelagiphage clones were obtained by using the dilution-to-  
457 extinction method (35). The purity of pelagiphages was verified by whole-genome  
458 sequencing.

459

460 **Morphological analysis by transmission electron microscopy.** Representative  
461 pelagiphages were observed by transmission electron microscopy (TEM).  
462 Pelagiphage lysates were filtered through 0.1µm pore-size filters and then  
463 concentrated using Amicon Ultra Centrifugal Filters (30-kDa, Merck Millipore).  
464 Concentrated phage particles were absorbed onto copper grids in the dark, negatively  
465 stained with 2% (wt/vol) uranyl acetate for two minutes, and air-dried. Samples were  
466 observed using a Hitachi transmission electron microscope at an acceleration voltage  
467 of 80 kV.

468

469 **Phage DNA preparation, genome sequencing and functional annotation.** Phage  
470 lysate preparation and concentration were carried out as described in Zhang and  
471 colleagues(35). Briefly, 250 ml of each phage lysate was filtered through 0.1 µm  
472 Supor membrane to remove cells and cell debris. Phage lysates were concentrated by  
473 centrifugal filtration using Amicon Ultra-15 30-kDa filters (Merck Millipore, Cork,  
474 Ireland). Phage genomic DNA was prepared using a formamide, phenol/chloroform  
475 extraction protocol(41) and sequenced on an Illumina paired-end HiSeq 2500  
476 platform. The raw reads were quality-filtered, trimmed and de novo assembled with  
477 default settings using CLC Genomic Workbench 11.0.1 (QIAGEN, Hilden, Germany).  
478 The remaining gaps in each pelagiphage genome were closed by Sanger sequencing  
479 of PCR products.

480 Prodigal(42) and GeneMark (43) were used for phage open reading frames  
481 (ORFs) prediction. The translated open reading frames (ORFs) were used as BLASTP  
482 queries to search against the NCBI nonredundant (nr) and NCBI-Refseq database.  
483 Putative functions were assigned to ORFs based on their homology to proteins of  
484 known function. In this study, genes with  $\geq 25\%$  amino acid identity,  $\geq 50\%$  alignment

485 coverage of the shortest protein, and an E-value cutoff  $\leq 1E-3$  were considered to be  
486 putative homologues. A PFAM database search was performed to identify conserved  
487 protein domains. HHPred was also employed to identify the distant protein homologs.  
488 tRNA prediction was performed using the tRNAscan-SE program (44). Comparative  
489 genome map and connections between homologous genes were visualized using  
490 CIRCOs(45).

491 The genomic sequences of the ten pelagiphages have been deposited in GenBank  
492 under the accession numbers MN698239 to MN698248.

493

494 **Network analysis.** Protein sequences from a total of 2591 bacterial viruses' genomes  
495 were downloaded from NCBI-RefSeq (v96). 927 viral genomic sequences ( $\geq 20$ kb,  
496  $\geq 20\%$  shared gene with any pelagiphage genome) from metagenomic fosmids, GOV  
497 and GOV2.0 datasets were also included in the network analysis(5, 6, 9). All proteins  
498 were compared using all-verses-all BLASTP (e-value  $\leq 1E-5$ , bitscore  $\geq 50$ ). Protein  
499 clusters (PCs) were defined using the Markov clustering algorithm (MCL) (46).  
500 vConTact 2.0 was then used to calculate a similarity score between every pair of  
501 genomes based on the number of PCs shared between two sequences and all pairs  
502 using the hypergeometric similarity (47). The network was created using Cytoscape  
503 v.3.5.1 (48).

504

505 **Phylogenetic analysis.** A phylogenetic tree of DNA polymerase family A domain  
506 sequences was constructed to evaluate the evolutionary relationship among  
507 pelagiphages and other diverse phages. Alignment for the DNA polymerase family A  
508 domain was constructed with MUSCLE (49) and edited with Gblock (50). The  
509 alignment was evaluated for optimal amino acid substitution models using ProtTest  
510 (51), and run with RAxML v8 (52) with a bootstrap of 500.

511 In order to evaluate the phylogenetic relationship among pelagiphages and  
512 environmental viral sequences, we constructed maximum likelihood phylogenetic  
513 trees of DNA polymerases and capsid proteins. Sequence alignments and editing were

514 performed using MUSCLE (49) and Gblocks (50), respectively. Maximum-likelihood  
515 phylogenetic trees were constructed using FastTree 2.1 (53) with WAG substitution  
516 model for amino acids.

517

518 **Search for integrase genes in pelagiphage related viral sequences.** Environmental  
519 viral genomic sequences from the viral clusters VC\_009, VC\_005, VC\_016 and  
520 VC\_018 that generated by vConTact 2.0 in this study were used for analysis. To  
521 identify integrase genes, hmmbuild was used to build HMM files from phage  
522 integrase and recombinase domains. The program hmmsearch was then used to  
523 identify the putative integrase genes by searching HMM files against the pelagiphage  
524 related environmental viral sequences. Viral sequences containing the putative  
525 integrase were subjected to manual inspection and comparative genomic analyse.  
526 Putative integration sites were identified by searching against known SAR11 genome  
527 sequences using BLASTn.

528

529 **Metagenomic fragment recruitment analyses.** Marine viromic datasets that were  
530 used for accessing phage relative abundance include Pacific Ocean Virome (POV),  
531 Scripps Pier Virome (SPV), India Ocean Virome (IOV), Malaspina Expedition  
532 virome (ME) and Global Oceans Viromes (GOV) (*S1 Appendix*, Table S1). The  
533 fragment reciprocal recruitment method was described in detail in a previous study  
534 (35). The analysis procedure is summarized as follows:

- 535 1. Each of the marine viromic reads was searched as a query against the NCBI-  
536 RefSeq viral database (release 88), 18 recently published HTVC019P-type  
537 pelagiphage and RCA phage genomes (22, 35), 6 newly sequenced HTVC010P-  
538 type pelagiphage genomes and 10 pelagiphage genomes reported in this study  
539 using DIAMOND BLASTx (e-value cutoff of  $\leq 1E-3$ , bitscore cutoff  $\geq 40$ ). Reads  
540 with homology to viral sequences were retained for the subsequent analysis.
- 541 2. The reads were assigned to the best-hit virus or best-hit bacteria using BLASTx  
542 against RefSeq viral database, RefSeq bacterial database, 11 HTVC019P-type

543 pelagiphages, 7 RCA phage genomes, 6 new HTVC010P-type pelagiphages and  
544 10 newly sequenced pelagiphages.

- 545 3. Reads that returned a best-hit of the query genome from the bacteriophage  
546 genomes included in the relative abundance analyses (*SI Appendix*, Table S2)  
547 were identified and extracted from the viromic datasets.
- 548 4. The relative abundances of each phage group were calculated and normalized as  
549 the number of reads recruited to the group normalized to the total number of base  
550 pairs in the virome and the average genome size (Reads mapped per kb per  
551 millions of reads, RPKM).

552 Due to the large amount of sequencing data in the Global Oceans Viromes  
553 (GOV) datasets (>500G), a different metagenomic analysis strategy was used to  
554 determine the relative abundances of different phage groups in the GOV. GOV reads  
555 were recruited onto the phage genomes (*SI Appendix*, Table S2) using BLASTx with  
556 an e-value cutoff of 1E-10. If a read was recruited to more than one phage genome,  
557 the read was associated with the phage that provided the highest bitscore. RPKM was  
558 also used to calculate and normalize the relative abundances of each phage group in  
559 each virome datasets in GOV.

560

561 **Data Availability Statement.** The genomic sequences of the ten pelagiphages have  
562 been deposited in GenBank under the accession numbers MN698239 to MN698248.  
563 The data that support the findings of this study are available upon request from the  
564 corresponding authors.

565

### 566 **Acknowledgments**

567 The work has been supported by NSFC grant 41706173. We thank Sijun Huang for  
568 providing the water samples. We thank Chen Li and Sun Jing for their assistance in  
569 TEM.

570

### 571 **Competing interests**

572 The authors declare that they have no conflict of interest.

573

574 **References**

- 575 1. K. E. Wommack, R. R. Colwell, Virioplankton: viruses in aquatic ecosystems.  
576 *Microbiol. Mol. Biol. Rev.* **64**, 69-114 (2000).
- 577 2. J. A. Fuhrman, Marine viruses and their biogeochemical and ecological effects.  
578 *Nature* **399**, 541-548 (1999).
- 579 3. C. A. Suttle, Marine viruses--major players in the global ecosystem. *Nat. Rev.*  
580 *Microbiol.* **5**, 801-812 (2007).
- 581 4. J. R. Brum *et al.*, Patterns and ecological drivers of ocean viral communities.  
582 *Science* **348**, 1261498 (2015).
- 583 5. S. Roux *et al.*, Ecogenomics and potential biogeochemical impacts of globally  
584 abundant ocean viruses. *Nature* **537**, 689-693 (2016).
- 585 6. A. C. Gregory *et al.*, Marine DNA viral macro- and microdiversity from pole to  
586 pole. *Cell* **177**, 1109-1123 (2019).
- 587 7. D. Paez-Espino *et al.*, Uncovering Earth's virome. *Nature* **536**, 425-430 (2016).
- 588 8. Y. Nishimura *et al.*, Environmental viral genomes shed new light on virus-host  
589 interactions in the ocean. *mSphere* **2**, e00359-16 (2017).
- 590 9. C. M. Mizuno, F. Rodriguezvalera, N. E. Kimes, R. Ghai, Expanding the marine  
591 virosphere using metagenomics. *PLoS Genet.* **9**, e1003987 (2013).
- 592 10. C. M. Mizuno, R. Ghai, A. Saghaï, P. Lopez-Garcia, F. Rodriguez-Valera,  
593 Genomes of abundant and widespread viruses from the deep ocean. *mBio* **7**,  
594 e00805-16 (2016).
- 595 11. S. Roux *et al.*, Ecology and evolution of viruses infecting uncultivated SUP05  
596 bacteria as revealed by single-cell- and meta-genomics. *eLife* **3**, e03125 (2014).
- 597 12. J. M. Labonte *et al.*, Single-cell genomics-based analysis of virus-host  
598 interactions in marine surface bacterioplankton. *ISME J.* **9**, 2386-2399 (2015).
- 599 13. F. Martinez-Hernandez *et al.*, Single-virus genomics reveals hidden cosmopolitan  
600 and abundant viruses. *Nat. Commun.* **8**, 15892 (2017).

- 601 14. J. Ren, N. A. Ahlgren, Y. Y. Lu, J. A. Fuhrman, F. Sun, VirFinder: a novel k-mer  
602 based tool for identifying viral sequences from assembled metagenomic data.  
603 *Microbiome* **5**, 69 (2017).
- 604 15. I. Kang, H. M. Oh, D. Kang, J. C. Cho, Genome of a SAR116 bacteriophage  
605 shows the prevalence of this phage type in the oceans. *Proc. Natl. Acad. Sci.*  
606 *U.S.A.* **110**, 12343-12348 (2013).
- 607 16. Y. Zhao *et al.*, Abundant SAR11 viruses in the ocean. *Nature* **494**, 357-360,  
608 (2013).
- 609 17. A. I. Culley, Insight into the unknown marine virus majority. *Proc. Natl. Acad.*  
610 *Sci. U.S.A.* **110**, 12166-12167 (2013).
- 611 18. M. S. Rappe, S. A. Connon, K. L. Vergin, S. J. Giovannoni, Cultivation of the  
612 ubiquitous SAR11 marine bacterioplankton clade. *Nature* **418**, 630-633, (2002).
- 613 19. S. J. Giovannoni, SAR11 Bacteria: The most abundant plankton in the oceans.  
614 *Ann. Rev. Mar. Sci.* **9**, 231-255 (2017).
- 615 20. D. Tsementzi *et al.*, SAR11 bacteria linked to ocean anoxia and nitrogen loss.  
616 *Nature* **536**, 179-183 (2016).
- 617 21. P. Carini, L. Steindler, S. Beszteri, S. J. Giovannoni, Nutrient requirements for  
618 growth of the extreme oligotroph '*Candidatus Pelagibacter ubique*' HTCC1062 on  
619 a defined medium. *ISME J.* **7**, 592-602 (2013).
- 620 22. Y. Zhao *et al.*, Pelagiphages in the *Podoviridae* family integrate into host  
621 genomes. *Environ. Microbiol.* **21**, 1989-2001 (2019).
- 622 23. O. Fayet, T. Ziegelhoffer, C. Georgopoulos, The *groES* and *groEL* heat shock  
623 gene products of *Escherichia coli* are essential for bacterial growth at all  
624 temperatures. *J. Bacteriol.* **171**, 1379-1385 (1989).
- 625 24. P. A. Lund, Multiple chaperonins in bacteria--why so many? *FEMS Microbiol.*  
626 *Rev.* **33**, 785-800 (2009).
- 627 25. S. M. van der Vies, A. A. Gatenby, C. Georgopoulos, Bacteriophage T4 encodes a  
628 co-chaperonin that can substitute for *Escherichia coli* GroES in protein folding.  
629 *Nature* **368**, 654-656 (1994).



- 630 26. D. Ang *et al.*, Pseudo-T-even bacteriophage RB49 encodes CocO, a cochaperonin  
631 for GroEL, which can substitute for *Escherichia coli's* GroES and bacteriophage  
632 T4's Gp31. *J. Biol. Chem.* **276**, 8720-8726 (2001).
- 633 27. R. L. Marine, D. J. Nasko, J. Wray, S. W. Polson, K. E. Wommack, Novel  
634 chaperonins are prevalent in the viroplankton and demonstrate links to viral  
635 biology and ecology. *ISME J.* **11**, 2479-2491 (2017).
- 636 28. H. F. Schmidt, E. G. Sakowski, S. J. Williamson, S. W. Polson, K. E. Wommack,  
637 Shotgun metagenomics indicates novel family A DNA polymerases predominate  
638 within marine viroplankton. *ISME. J.* **8**, 103-114 (2014).
- 639 29. F. Martinez-Hernandez *et al.*, Single-cell genomics uncover Pelagibacter as the  
640 putative host of the extremely abundant uncultured 37-F6 viral population in the  
641 ocean. *ISME J.* **13**, 232-236. (2018).
- 642 30. D. Esposito, J. J. Scocca, The integrase family of tyrosine recombinases:  
643 evolution of a conserved active site domain. *Nucleic Acids Res.* **25**, 3605-3614  
644 (1997).
- 645 31. S. E. Nunes-Düby, H. J. Kwon, R. S. Tirumalai, T. Ellenberger, A. Landy,  
646 Similarities and differences among 105 members of the Int family of site-specific  
647 recombinases. *Nucleic Acids Res.* **26**, 391-406 (1998).
- 648 32. B. L. Hurwitz, M. B. Sullivan, The Pacific Ocean virome (POV): a marine viral  
649 metagenomic dataset and associated protein clusters for quantitative viral  
650 ecology. *PLoS One* **8**, e57355 (2013).
- 651 33. V. Bischoff *et al.*, Cobaviruses - a new globally distributed phage group infecting  
652 *Rhodobacteraceae* in marine ecosystems. *ISME J.* **13**, 1404-1421 (2019).
- 653 34. N. A. Ahlgren, C. A. Fuchsman, G. Rocap, J. A. Fuhrman, Discovery of several  
654 novel, widespread, and ecologically distinct marine *Thaumarchaeota* viruses that  
655 encode *amoC* nitrification genes. *ISME J.* **13**, 618-631, (2019).
- 656 35. Z. Zhang *et al.*, Diverse, abundant and novel viruses infecting “unculturable” but  
657 abundant marine bacteria. Preprint *bioRxiv*  
658 <https://www.biorxiv.org/content/10.1101/699256v1> (2019).

- 659 36. S. A. Connon, S. J. Giovannoni, High-throughput methods for culturing  
660 microorganisms in very-low-nutrient media yield diverse new marine isolates.  
661 *Appl. Environ. Microbiol.* **68**, 3878-3885 (2002).
- 662 37. R. R. Malmstrom, R. P. Kiene, M. T. Cottrell, D. L. Kirchman, Contribution of  
663 SAR11 bacteria to dissolved dimethylsulfoniopropionate and amino acid uptake  
664 in the North Atlantic ocean. *Appl. Environ. Microbiol.* **70**, 4129-4135 (2004).
- 665 38. J. C. Thrash *et al.*, Single-cell enabled comparative genomics of a deep ocean  
666 SAR11 bathytype. *ISME J.* **8**, 1440-1451 (2014).
- 667 39. K. L. Vergin *et al.*, High-resolution SAR11 ecotype dynamics at the Bermuda  
668 Atlantic Time-series Study site by phylogenetic placement of pyrosequences.  
669 *ISME J.* **7**, 1322-1332 (2013).
- 670 40. I. Kang, J. C. Cho, Depth-specific distribution of the SAR116 phages revealed by  
671 virome binning. *J. Microbiol. Biotechnol.* **24**, 592-596 (2014).
- 672 41. J. Sambrook, E. F. Fritsch, T. Maniatis, *Molecular Cloning: A Laboratory*  
673 *Manual* (Cold Spring Harbor Laboratory, Cold Spring Harbor, 1989).
- 674 42. D. Hyatt *et al.*, Prodigal: prokaryotic gene recognition and translation initiation  
675 site identification. *BMC Bioinform.* **11**, 119 (2010).
- 676 43. M. Borodovsky, J. McIninch, GENMARK: Parallel gene recognition for both  
677 DNA strands. *Comput. Chem.* **17**, 123-133 (1993).
- 678 44. T. M. Lowe, S. R. Eddy, tRNAscan-SE: a program for improved detection of  
679 transfer RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955-964 (1997).
- 680 45. M. Krzywinski *et al.*, Circos: an information aesthetic for comparative genomics.  
681 *Genome Res.* **19**, 1639-1645 (2009).
- 682 46. A. J. Enright, S. Van Dongen, C. A. Ouzounis, An efficient algorithm for large-  
683 scale detection of protein families. *Nucleic Acids Res.* **30**, 1575-1584 (2002).
- 684 47. H. Bin Jang *et al.*, Taxonomic assignment of uncultivated prokaryotic virus  
685 genomes is enabled by gene-sharing networks. *Nat. Biotechnol.* **37**, 632-639  
686 (2019).
- 687 48. P. Shannon, *et al.*, Cytoscape: a software environment for integrated models of

- 688        biomolecular interaction networks. *Genome Res.* **13**, 2498-2504 (2003).
- 689    49. R. C. Edgar, MUSCLE: multiple sequence alignment with high accuracy and  
690        high throughput. *Nucleic Acids Res.* **32**, 1792-1797 (2004).
- 691    50. J. Castresana, Selection of conserved blocks from multiple alignments for their  
692        use in phylogenetic analysis. *Mol. Biol. Evol.* **17**, 540-552 (2000).
- 693    51. F. Abascal, R. Zardoya, D. Posada, ProtTest: selection of best-fit models of  
694        protein evolution. *Bioinformatics* **21**, 2104-2105 (2005).
- 695    52. A. Stamatakis, RAxML version 8: a tool for phylogenetic analysis and post-  
696        analysis of large phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).
- 697    53. M. N. Price, P. S. Dehal, A. P. Arkin, FastTree 2--approximately maximum-  
698        likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).

699 **Figure legends**

700 **Fig. 1.** (A) Transmission electron microscope images of selected pelagiphages from  
701 each group. A, HTVC023P; B, HTVC027P; C, HTVC103P; D, HTVC104P; E,  
702 HTVC111P; F, HTVC106P. (Scale bars: 50 nm). (B) Heatmap presentation of shared  
703 genes of newly isolated pelagiphages. Phages in the same group are boxed. (C)  
704 Unrooted maximum-likelihood phylogenetic tree of phage DNA polymerases  
705 constructed with conserved polymerase domains. Novel DNA polymerases identified  
706 from a previous study are in bold (28). The scale bar represents the amino acid  
707 substitutions per site.

708 **Fig. 2.** Genomic organization and functional annotation of distinct pelagiphage  
709 genera. (A) HTVC023P-type pelagiphage genomes, (B) HTVC103P-type  
710 pelagiphages are compared to SAR116 phage HMO-2011, (C) HTVC111P-type  
711 pelagiphage genomes, (D) HTVC106P genome. Predicted ORFs are indicated by  
712 arrows and color-coded according to their putative biological function. Homologous  
713 genes were connected by dashed lines. The color of the shading connecting  
714 homologous genes indicate the level of amino acid identities between genes. The  
715 arrows also designates the direction of transcription. Abbreviation: RNAP, RNA  
716 polymerase; SSB, single-stranded DNA binding protein; endo, endonuclease; DNAP,  
717 DNAP polymerase; exon, exonuclease; MazG, pyrophosphatase; DNA cytosine  
718 methyltransferase; FkbM family methyltransferase; TerS, terminase small subunit;  
719 TerL, terminase large subunit; GroEs, Co-Chaperonin GroES; HlyC, toxin-activating  
720 lysine-acyltransferase; purM, phosphoribosylaminoimidazole synthetase; GTF,  
721 Glycosyltransferase.

722 **Fig. 3.** Gene-content-based viral network of pelagiphages, related bacteriophages  
723 from NCBI, and related environmental viral sequences from Mediterranean DCM  
724 (MedDCM) fosmids, GOV and GOV2.0. The nodes represent the viral genomic  
725 sequences. The edges represent the similarities score between genomes based on  
726 shared gene content. Viral genomes that belong to different viral clusters are indicated  
727 by different colors. For clarity, only environmental viral sequences grouped with ten

728 pelagiphages were presented, and only bacteriophages have genome-genome  
729 similarity score of  $\geq 1$  with these four phage clusters were presented. Viral clusters  
730 generated by vConTACT2 are provided in Dataset S2 in the supplemental material.

731 **Fig. 4.** Maximum-likelihood tree of DNA polymerases from the viral cluster VC\_009  
732 generated by vConTACT v.2.0 in this study. HTVC023P-type pelagiphages, GOV  
733 populations (VC\_6 and VC\_8) and outgroups are indicated in green, red, and blue,  
734 respectively. The names of the GOV2.0 populations are omitted in the tree.

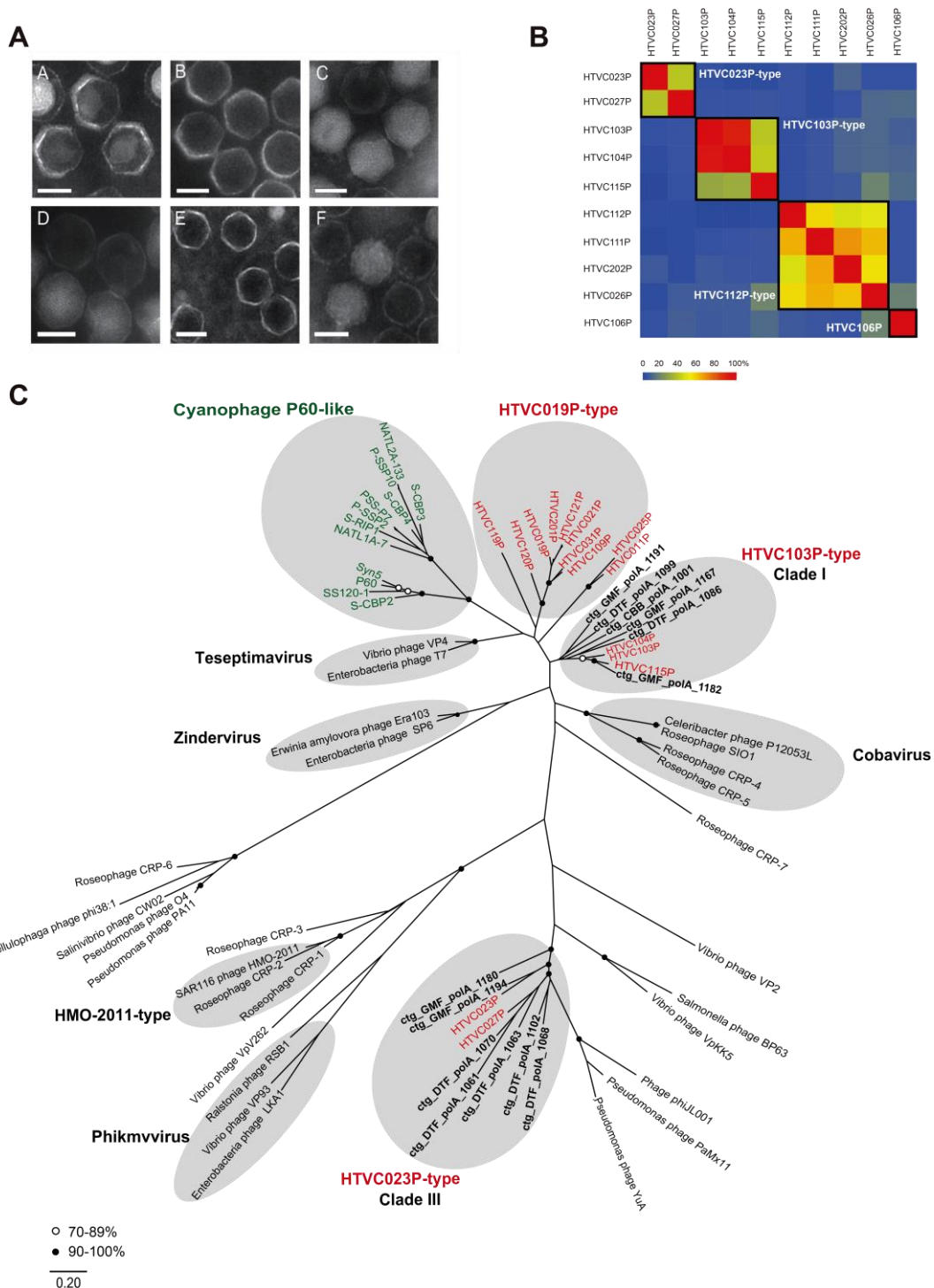
735 **Fig. 5.** Circos comparison plot indicating the genome comparison between  
736 pelagiphages and environmental viral fragments. (A) Comparison of HTVC023P-type  
737 genomes, vSAG 37-F6, and representative viral populations in VC\_6 and VC\_8. (B)  
738 Comparison of HTVC111P-type genomes and representative viral populations in  
739 VC\_41. (C) Comparison of HTVC106P genome and representative viral populations  
740 in VC\_67. Each coloured segment represents a phage genome or viral fragment with  
741 the numbers on the external surface indicating genome size in kb. Homologous genes  
742 shared between genomes are connected by color lines. Only the relatedness between  
743 pelagiphages and other sequences is indicated.

744 **Fig. 6.** Box plots indicate the relative abundance of major phage groups in different  
745 marine viromic datasets. Normalized read recruitment is depicted as the number of  
746 reads recruited per kilobase of the genome per billions reads in the dataset.

747 Pelagiphage represented groups are colored in blue; the HMO-2011-type group is  
748 colored in red. Pelagiphage groups identified in this study are marked with blue stars.

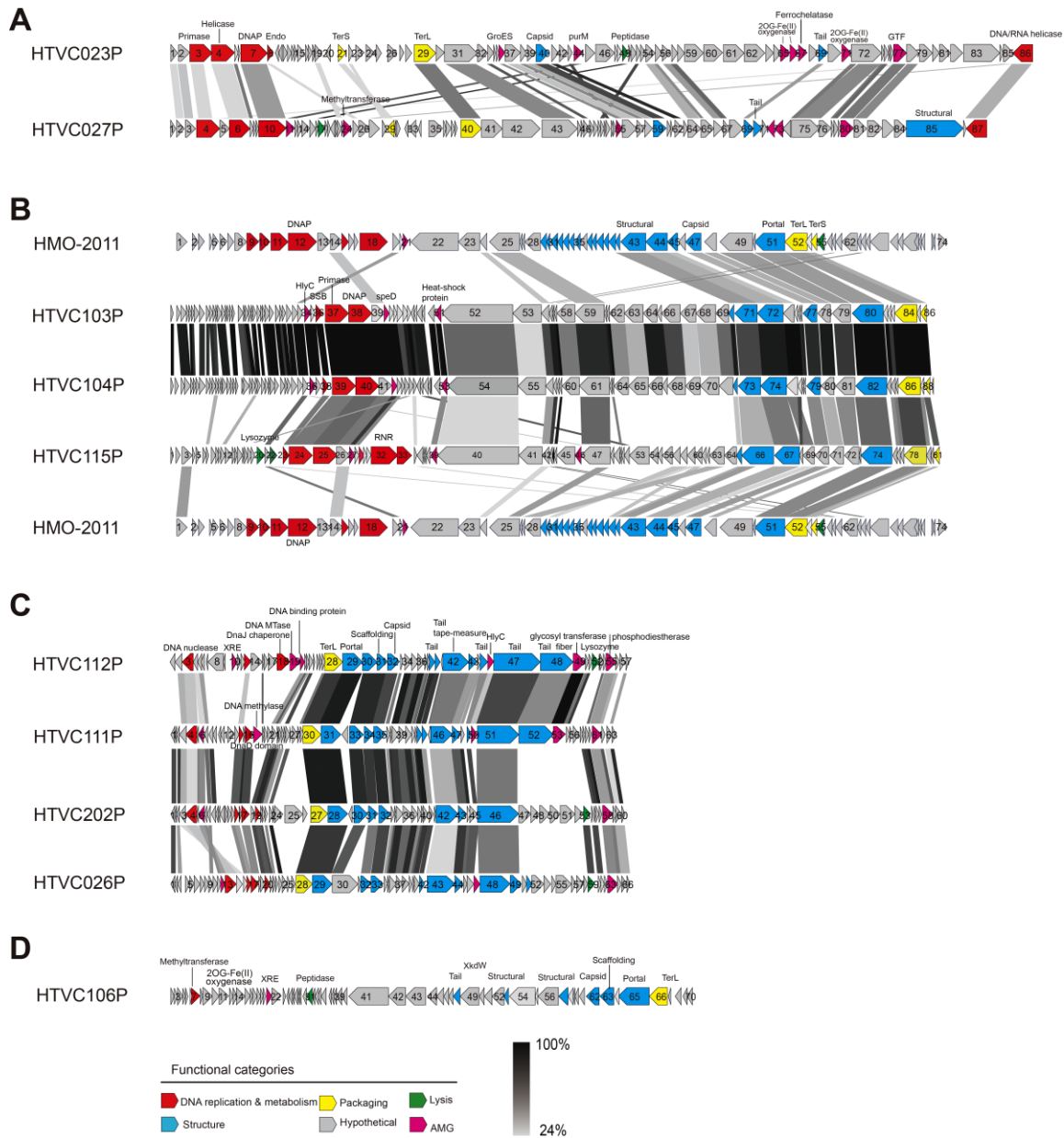
749 (A) Relative abundance of major phage groups in epipelagic, mesopelagic and  
750 bathypelagic samples in POV, MES, SPV and IOV. (B) Relative abundance of major  
751 phage groups in epipelagic, mesopelagic and bathypelagic samples in the Global  
752 Oceans Viromes (GOV). EPI, Epipelagic; MES, Mesopelagic; BAT, Bathypelagic.

753 **Fig. 7.** Comparison of the relative abundance between HTVC023P-type and HMO-  
754 2011-type, HTVC023P-type and HTVC010P-type. Upper panel: comparison of the  
755 relative abundance in POV, MES, SPV and IOV. Lower panel: comparison of the  
756 relative abundance in GOV. EPI, Epipelagic; MES, Mesopelagic; BAT, Bathypelagic



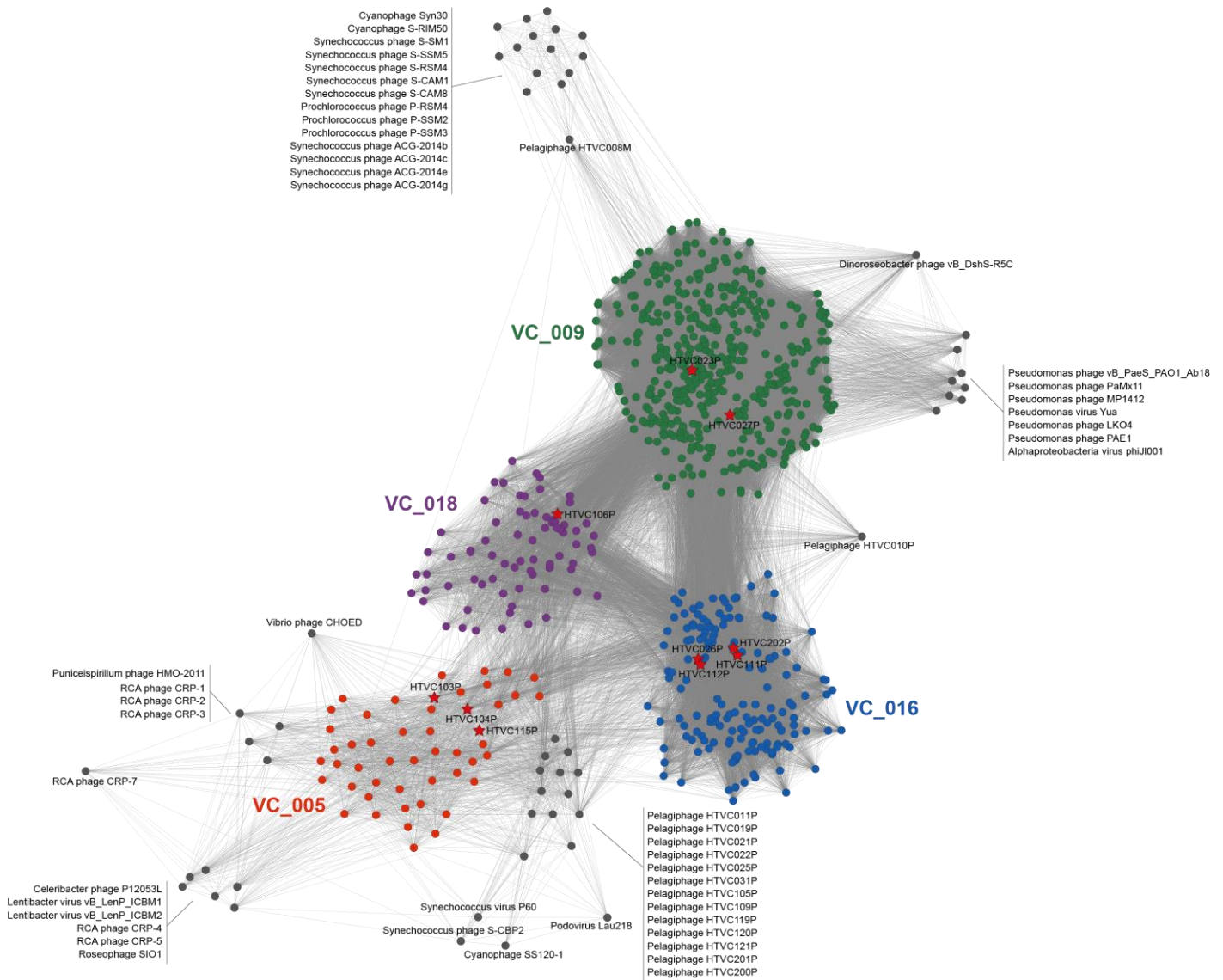
**Fig. 1.** (A) Transmission electron microscope images of selected pelagiphages from each group. A, HTVC023P; B, HTVC027P; C, HTVC103P; D, HTVC104P; E, HTVC111P; F, HTVC106P. (Scale bars: 50 nm). (B) Heatmap presentation of shared genes of newly isolated pelagiphages. Phages in the same group are boxed. (C) Unrooted maximum-likelihood phylogenetic tree of phage DNA polymerases constructed with conserved polymerase domains. Novel DNA polymerases identified from a previous study are in bold (28). The scale bar represents the amino acid substitutions per site.



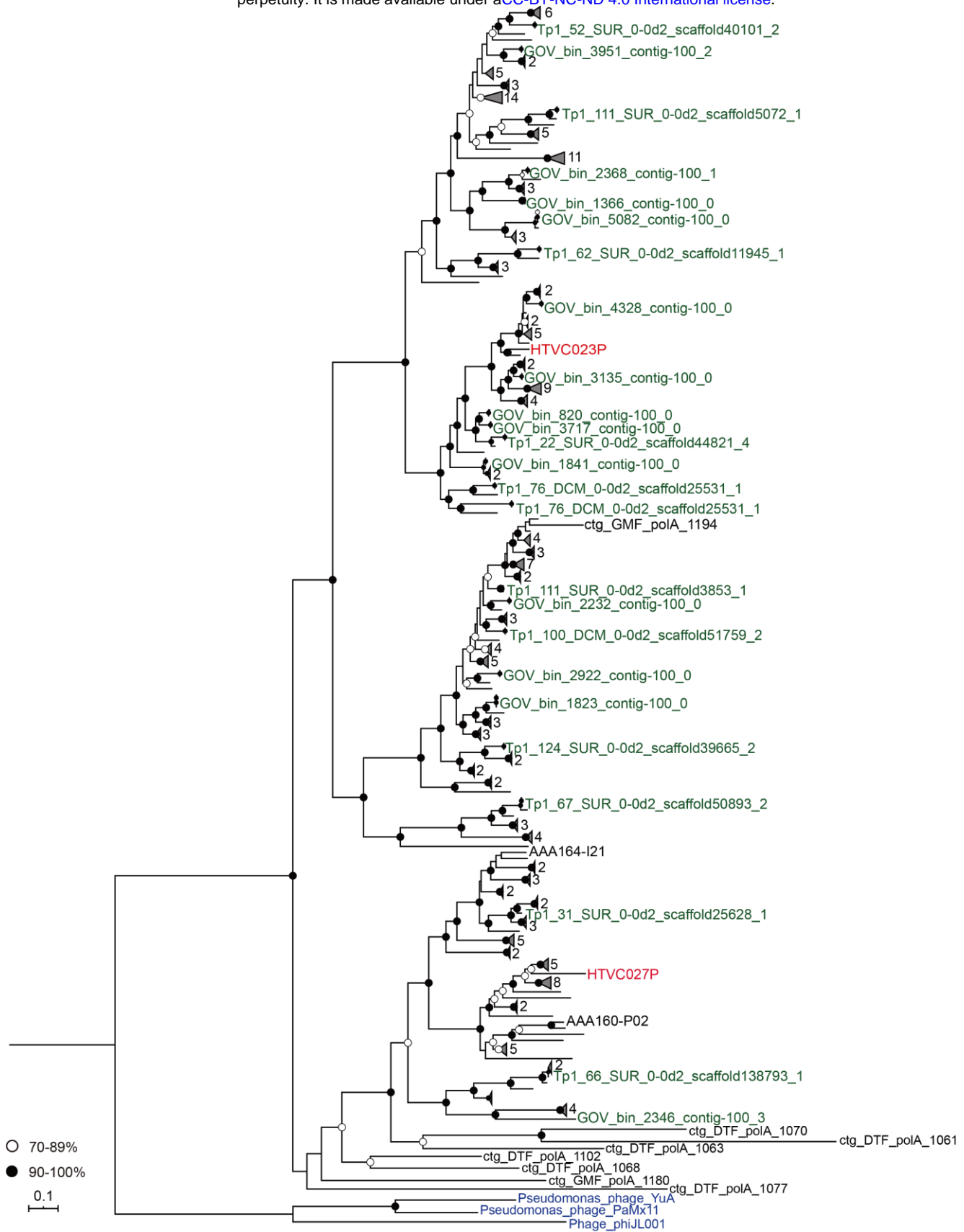


**Fig. 2.** Genomic organization and functional annotation of distinct pelagiphage genera. (A) HTVC023P-type pelagiphage genomes, (B) HTVC103P-type pelagiphages are compared to SAR116 phage HMO-2011, (C) HTVC111P-type pelagiphage genomes, (D) HTVC106P genome. Predicted ORFs are indicated by arrows and color-coded according to their putative biological function. Homologous genes were connected by dashed lines. The color of the shading connecting homologous genes indicate the level of amino acid identities between genes. The arrows also designates the direction of transcription. Abbreviation: RNAP, RNA polymerase; SSB, single-stranded DNA binding protein; endo, endonuclease; DNAP, DNAP polymerase; exon, exonuclease; MazG, pyrophosphatase; DNA cytosine methyltransferase; FkbM family methyltransferase; TerS, terminase small subunit; TerL, terminase large subunit; GroEs, Co-Chaperonin GroES; HlyC, toxin-activating lysine-acyltransferase; purM, phosphoribosylaminoimidazole synthetase; GTF, Glycosyltransferase.

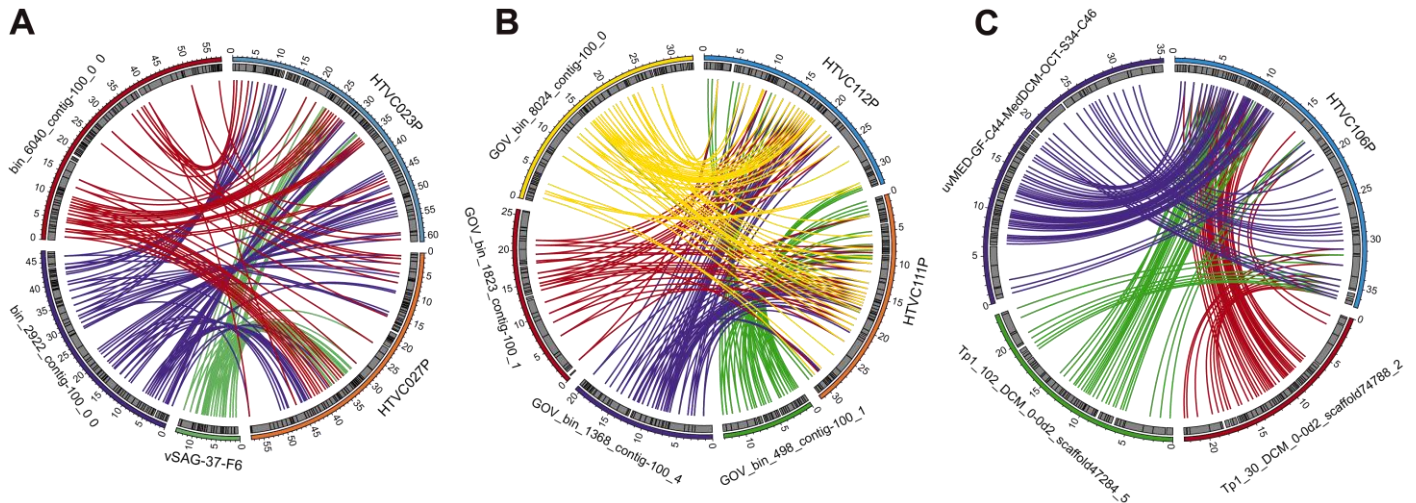




776  
777  
778 **Fig. 3.** Gene-content-based viral network of pelagiphages, related bacteriophages from NCBI, and related  
779 environmental viral sequences from Mediterranean DCM (MedDCM) fosmids, GOV and GOV2.0. The  
780 nodes represent the viral genomic sequences. The edges represent the similarities score between genomes  
781 based on shared gene content. Viral genomes that belong to different viral clusters are indicated by different  
782 colors. For clarity, only environmental viral sequences grouped with ten pelagiphages were presented, and  
783 only bacteriophages have genome-genome similarity score of  $\geq 1$  with these four phage clusters were  
784 presented. Viral clusters generated by vConTACT2 are provided in Dataset S2 in the supplemental material.



785  
786 **Fig. 4.** Maximum-likelihood tree of DNA polymerases from the viral cluster VC\_009 generated by  
787 vConTACT v.2.0 in this study. HTVC023P-type pelagiphages, GOV populations (VC\_6 and VC\_8) and  
788 outgroups are indicated in green, red, and blue, respectively. The names of the GOV2.0 populations are  
789 omitted in the tree.



790

791

792

793

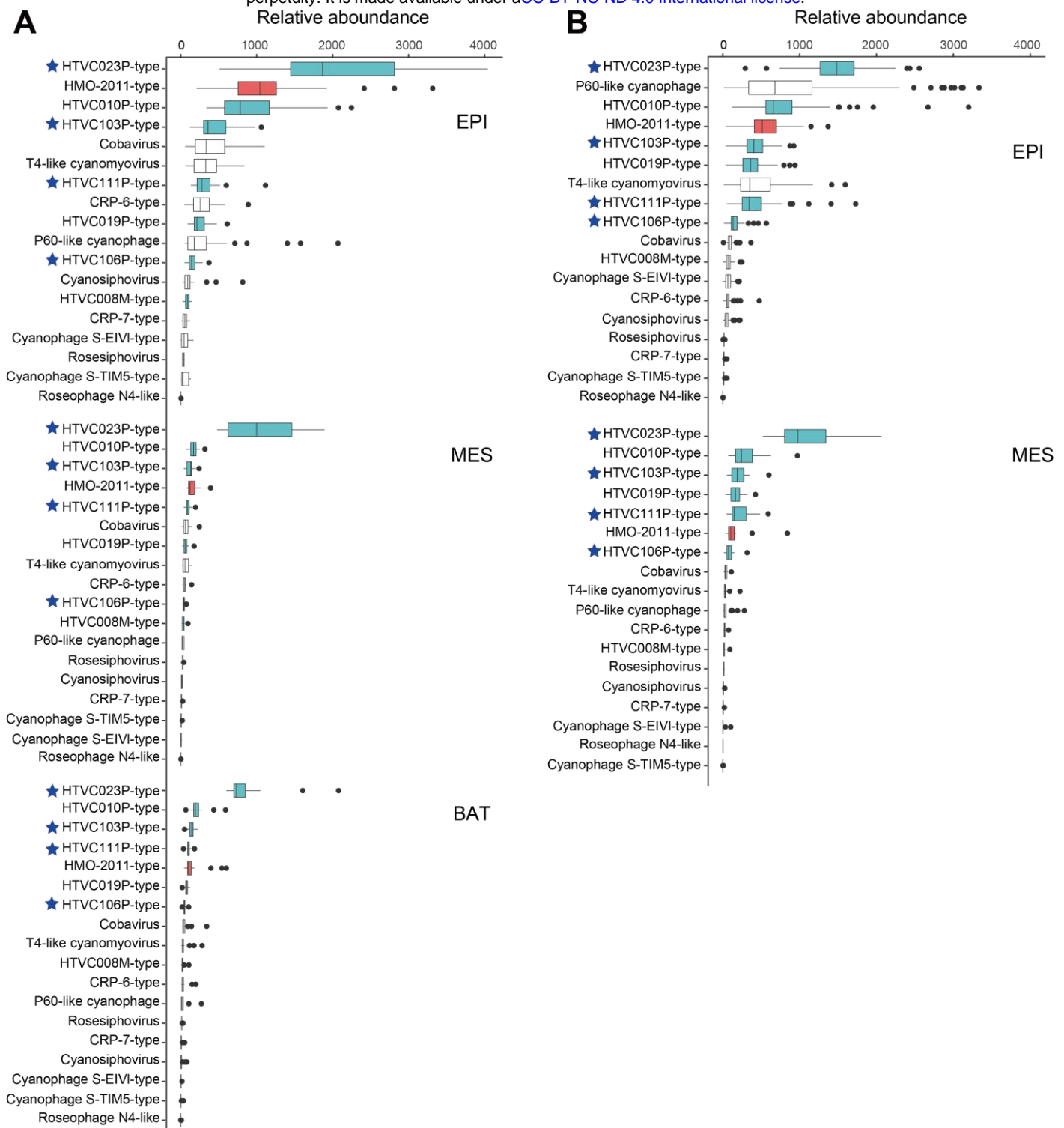
794

795

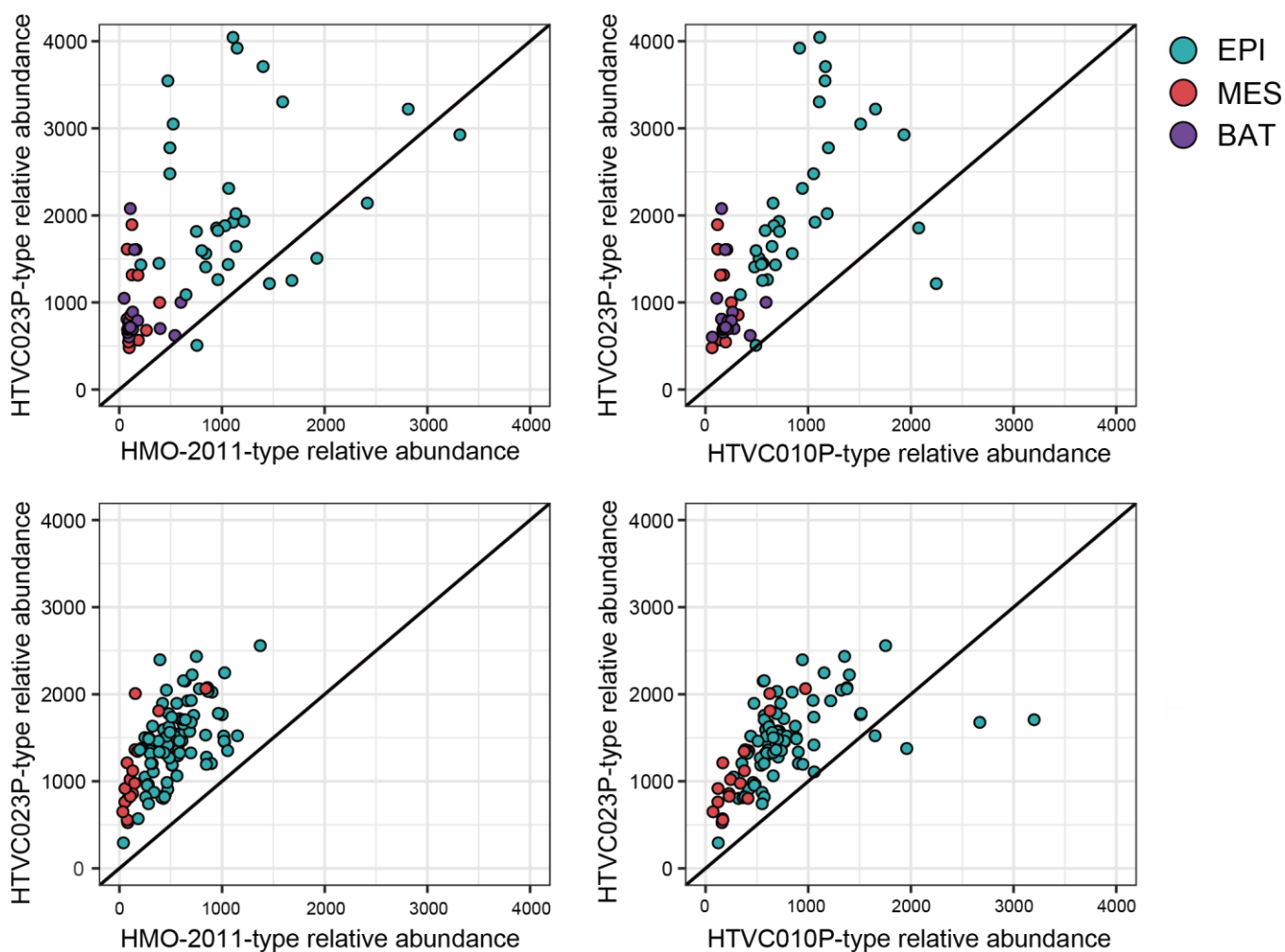
796

797

**Fig. 5.** Circos comparison plot indicating the genome comparison between pelagiphages and environmental viral fragments. (A) Comparison of HTVC023P-type genomes, vSAG 37-F6, and representative viral populations in VC\_6 and VC\_8. (B) Comparison of HTVC111P-type genomes and representative viral populations in VC\_41. (C) Comparison of HTVC106P genome and representative viral populations in VC\_67. Each coloured segment represents a phage genome or viral fragment with the numbers on the external surface indicating genome size in kb. Homologous genes shared between genomes are connected by color lines. Only the relatedness between pelagiphages and other sequences is indicated.



**Fig. 6.** Box plots indicate the relative abundance of major phage groups in different marine viromic datasets. Normalized read recruitment is depicted as the number of reads recruited per kilobase of the genome per billions reads in the dataset. Pelagiphage represented groups are colored in blue; the HMO-2011-type group is colored in red. Pelagiphage groups identified in this study are marked with blue stars. (A) Relative abundance of major phage groups in epipelagic, mesopelagic and bathypelagic samples in POV, MES, SPV and IOV. (B) Relative abundance of major phage groups in epipelagic, mesopelagic and bathypelagic samples in the Global Oceans Viromes (GOV). EPI, Epipelagic; MES, Mesopelagic; BAT, Bathypelagic.



806

807

808

809

810

**Fig. 7.** Comparison of the relative abundance between HTVC023P-type and HMO-2011-type, HTVC023P-type and HTVC010P-type. Upper panel: comparison of the relative abundance in POV, MES, SPV and IOV. Lower panel: comparison of the relative abundance in GOV. EPI, Epipelagic; MES, Mesopelagic; BAT, Bathypelagic.

811 **Table 1** General features of pelagiphages analyzed in this study.

Phage	Host	Source water	Depth (m)	Latitude	Longitude	Sampling date	Phage group	Capsid size (mean±s.d., nm)	Genome size (bp)	G+C %	Number of ORFs	Accession number
HTVC023P	HTCC1062	South China Sea SEATS	500	S18°00'	E116°00'	Nov-2016	HTVC023P-type	68±2	60878	35.0	86	MN698239
HTVC027P	HTCC1062	South Pole K1	150	N11°39'	E78°59'	Sep-2016	HTVC023P-type	69±1	57595	34.8	86	MN698241
HTVC103P	HTCC7211	South China Sea SEATS	5	S18°00'	E116°00'	Aug-2014	HTVC103P-type	68±1	54103	31.0	86	MN698242
HTVC104P	HTCC7211	India Ocean 105	75	S4°0'	E95°28'	Mar-2015	HTVC103P-type	67±2	54359	30.9	89	MN698243
HTVC115P	HTCC7211	India Ocean 105	500	S4°0'	E95°28'	Mar-2015	HTVC103P-type	69±3	54819	33.2	81	MN698247
HTVC111P	HTCC7211	Mediterranean Sea	Surface	N43°42'	W7°17'	Aug-2016	HTVC111P-type	55±2	31577	30.5	64	MN698245
HTVC112P	HTCC7211	South Pole K1	150	S78°59'	E11°40'	Sep-2016	HTVC111P-type	Nd	32478	30.4	58	MN698246
HTVC026P	HTCC1062	Mediterranean Sea	Surface	N43°42'	W7°17'	Aug-2016	HTVC111P-type	Nd	32480	31.3	66	MN698240
HTVC202P	FZCC0015	Pingtang coast, Taiwan straight	Surface	N25°26'	E119°47'	May-2017	HTVC111P-type	Nd	32226	31.3	61	MN698248
HTVC106P	HTCC7211	India Ocean 113 station	500	N0°0'	E92°59'	Mar-2015	HTVC106P-type	56±1	36945	32.1	70	MN698244

812 ND : Not Determined.

813