



HAL
open science

Assessing the viral content of uncultured picoeukaryotes in the global-ocean by single cell genomics

Yaiza M. Castillo, Jean-François Mangot, Luiz Felipe Benites, Ramiro Logares, Megumi Kuronishi, Hiroyuki Ogata, Olivier Jaillon, Ramon Massana, Marta Sebastian, Dolors Vaque

► To cite this version:

Yaiza M. Castillo, Jean-François Mangot, Luiz Felipe Benites, Ramiro Logares, Megumi Kuronishi, et al.. Assessing the viral content of uncultured picoeukaryotes in the global-ocean by single cell genomics. *Molecular Ecology*, 2019, 28 (18), pp.4272-4289. 10.1111/mec.15210 . cea-04307624

HAL Id: cea-04307624

<https://cea.hal.science/cea-04307624v1>

Submitted on 6 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Assessing the viral content of uncultured picoeukaryotes in the
2 global-ocean by single cell genomics
3

4 Yaiza M. Castillo^{1,*}, Jean-François Mangot^{1,*}, L. Felipe Benites², Ramiro Logares¹, Megumi
5 Kuronishi³, Hiroyuki Ogata³, Olivier Jaillon⁴, Ramon Massana¹, Marta Sebastián^{1,5} and
6 Dolors Vaqué¹

7
8 **Author affiliations**

9 ¹ Department of Marine Biology and Oceanography, Institute of Marine Sciences (ICM),
10 CSIC, Barcelona, Spain.

11 ² Integrative Biology of Marine Organisms (BIOM), Sorbonne University, CNRS,
12 Oceanological Observatory of Banyuls, Banyuls-sur-Mer, France.

13 ³ Bioinformatic Center, Institute for Chemical Research, Kyoto University, Gokasho, Uji,
14 Japan.

15 ⁴ Génomique Métabolique, Genoscope, Institut de biologie François Jacob, CEA, CNRS,
16 Université d'Evry, Université Paris-Saclay, Evry, France.

17 ⁵ Institute of Oceanography and Global Change (IOCAG), University of Las Palmas de Gran
18 Canaria, Telde, Spain.

19 ** These authors contributed equally to this work and thus joint first authors.*

20 **Correspondence:** Yaiza M. Castillo and Dolors Vaqué. Department of Marine Biology and
21 Oceanography, Institut of Marine Sciences (CSIC), Passeig Marítim de la Barceloneta, 37-49,
22 E-08003 Barcelona, Catalonia, Spain. E-mail: yaiza@icm.csic.es and dolors@icm.csic.es.

23 **Running title:** Viral signals in marine protists genomes

24 **Abstract**

25 Viruses are the most abundant biological entities on Earth and have fundamental ecological
26 roles in controlling microbial communities. Yet, although their diversity is being increasingly
27 explored, little is known about the extent of viral interactions with their protist hosts since
28 most studies are limited to a few cultivated species. Here, we exploit the potential of single-
29 cell genomics to unveil viral associations in 65 individual cells of 11 prevalent uncultured
30 stramenopiles lineages sampled during the *Tara* Oceans expedition. We identified viral
31 signals in 57% of the cells covering nearly every lineage and with narrow host specificity
32 signal. Only 7 out of the 64 detected viruses displayed homologies to known viral sequences.
33 A search for our viral sequences in global ocean metagenomes showed that they were
34 preferentially found at the DCM and within the 0.2-3 μm size fraction. Some of the viral
35 signals were widely distributed, while others geographically constrained. Among the viral
36 signals we detected an endogenous mavirus virophage potentially integrated within the
37 nuclear genome of two distant uncultured stramenopiles. Virophages have been previously
38 reported as a cell's defense mechanism against other viruses, and may therefore play an
39 important ecological role in regulating protist populations. Our results point to single-cell
40 genomics as a powerful tool to investigate viral associations in uncultured protists, suggesting
41 a wide distribution of these relationships, and providing new insights into the global viral
42 diversity.

43

44 **Keywords:** Single-cell genomics; viral associations; protists; uncultured stramenopiles;
45 viruses; virophages.

46 **Introduction**

47 Viruses are major players in marine biogeochemical cycles (Jover, Effler, Buchan, Wilhelm,
48 & Weitz, 2014) and constitute the most abundant biological entities in the oceans, ranging
49 from about 10^4 to 10^7 ml⁻¹ (Danovaro et al., 2011; Suttle, 2005). They are known to be a
50 major cause of microbial (bacteria, archaea and protists) mortalities (Munn, 2006), leading to
51 approximately 10^{29} infection events every day (Brussaard et al., 2008) and causing the
52 release of 10^8 - 10^9 tons of biogenic carbon per day (Brussaard et al., 2008; Suttle, 2005).
53 Furthermore, they are main vectors of gene transfer in the oceans (Middelboe & Brussaard,
54 2017), impacting microbial community dynamics, diversity and evolution (Breitbart, 2012;
55 Jover et al., 2014; Weitz & Wilhelm, 2012).

56 Our knowledge of marine viral diversity and biogeography has been constantly expanding
57 during this last decade with the advent of viral metagenomics (e.g., Coutinho et al., 2017;
58 Mizuno, Rodriguez-Valera, Kimes, & Ghai, 2013; Paez-Espino et al., 2016). Multiple studies
59 unveiled a large novel diversity of uncultured viruses, indicating their key roles in nutrient
60 cycling and trophic networks (Brum et al., 2015; Roux et al., 2016). Unfortunately, despite
61 these fruitful advances in viral ecology, our understanding of virus-host interactions is still in
62 its infancy. The question of ‘who infects whom’ within marine microbial communities has
63 always been central, and the assessment of the true extent of host specificity among marine
64 viruses remains challenging (Brum & Sullivan, 2015). For a long time, studies investigating
65 virus-host interactions were limited to cultured host cells, restricting our knowledge to the
66 0.1–1% of host cells that are in culture (Rappé & Giovannoni, 2003; Swan et al., 2013), and
67 biasing our knowledge towards virulent lytic viruses (Brüssow & Hendrix, 2002; Swan et al.,

68 2013). Thus, many viruses are still uncharacterized and novel culture-independent approaches
69 are needed to overcome these methodological limitations.

70 Several methods have been developed to investigate putative interactions between viruses and
71 uncultured hosts reviewed by Brum & Sullivan (2015) and Breitbart, Bonnain, Malki, &
72 Sawaya (2018). These include analyses by metaviromics (Bolduc, Wirth, Mazurie, & Young,
73 2015; Brum et al., 2015), matching CRISPR spacers (Anderson, Brazelton, & Baross, 2011;
74 Berg Miller et al., 2012), phageFISH (Fluorescence *In Situ* Hybridization; Allers et al., 2013),
75 viral tagging (Deng et al., 2012, 2014), the polony method (Baran, Goldin, Maidanik, &
76 Lindell, 2018), the use of microfluidic digital PCR (Tadmor, Ottesen, Leadbetter, & Phillips,
77 2011) and single-cell genomics (SCG) (e.g., Labonté et al., 2015 and Roux et al., 2014). From
78 these, SCG emerged as a powerful complement to cultivation and metagenomics by providing
79 genomic information from individual uncultured cells (Stepanauskas, 2012). Furthermore, it
80 has an incredible potential for cell-specific analyses of organismal interactions, such as
81 parasitism, symbiosis and predation (Krabberød, Bjorbækmo, Shalchian-Tabrizi, & Logares,
82 2017; Stepanauskas, 2012), giving comprehensive insights of *in situ* virus-host associations.
83 Indeed, this effective approach has revealed new associations between viruses and bacterial
84 (Labonté et al., 2015; Roux et al., 2014) or archaeal cells (Chow, Winget, White, Hallam, &
85 Suttle, 2015; Labonté et al., 2015; Munson-McGee et al., 2018). However, the application of
86 SCG to protist cells is relatively recent and there is still a limited number of Single-cell
87 Amplified Genomes (SAGs) from microeukaryotes (e.g., Bhattacharya et al., 2012; Heywood,
88 Sieracki, Bellows, Poulton, & Stepanauskas, 2011; Mangot et al., 2017; Roy et al., 2015;
89 Troell et al., 2016; Vannier et al., 2016), with only one study that has so far explored virus-
90 host interactions (Yoon et al., 2011).

91 In the present work, we use SCG to uncover putative interactions between viruses and
92 uncultured protists using 65 SAGs produced during the *Tara* Oceans expedition (Karsenti et
93 al., 2011). These cells were affiliated to 11 stramenopile lineages belonging to MARine
94 STramenopiles (MASTs), Chrysophyceae, Dictyochophyceae and Pelagophyceae, that are
95 known to be important components of marine pico- and nanosized eukaryotic assemblages (1-
96 5 μm , Massana, 2011). Initially detected in molecular diversity surveys, MASTs are formed
97 by at least 18 independent groups of essentially uncultured protists (Massana, del Campo,
98 Sieracki, Audic, & Logares, 2014), some of which display a widespread distribution in
99 sequencing data sets (Lin et al., 2012; Logares et al., 2012; Seeleuthner et al., 2018) and are
100 abundant in microscopy counts (Massana, Terrado, Forn, Lovejoy, & Pedrós-Alió,
101 2006). Within these MASTs, we analyzed SAGs from three clades, MAST-3, MAST-4 and
102 MAST-7 (Massana et al., 2014). We also report the putative linkages between viruses and
103 SAGs from the uncultured chrysophyte lineages G and H, formed by pigmented and colorless
104 cells respectively, which are abundant in molecular diversity surveys (del Campo & Massana,
105 2011; Seeleuthner et al., 2018). Finally, we screened for viral signatures in SAGs from a
106 cultured pelagophyte (*Pelagomonas calceolata*) and an uncultured dictyochophyte within the
107 order Pedinellales.

108 Our results revealed a large diversity of viral sequences associated to protist cells, the vast
109 majority of which correspond to previously unidentified viral lineages. Using global ocean
110 metagenomes from the *TARA* Oceans expedition, we looked at the geographical distribution
111 of the identified viral sequences in epipelagic waters by fragment recruitment analysis,
112 finding that some SAG-associated viruses were widely distributed while others are restricted
113 to certain areas. Finally, special attention was paid to a particular virophage sequence
114 retrieved in two distinct stramenopile lineages that is highly similar to the endogenous

115 *Cafeteriavirus*-dependent mavirus, known to be integrated within the nuclear genome of their
116 host (Fischer & Hackl, 2016). Overall, our approach constitutes an initial attempt to determine
117 virus-host associations within protists using culture-independent single-cell genomics.

118

119 **Materials and methods**

120 *Sample collection and single cell sorting*

121 Samples for single-cell sorting were collected during the circumglobal *Tara* Oceans
122 expedition (2009-2013) (Karsenti et al., 2011) and processed as described in Alberti et al.,
123 2017. Flow cytometry cell sorting on cryopreserved samples and genomic DNA amplification
124 by multiple displacement amplification (MDA) were performed at the Single Cell Genomics
125 Center in the Bigelow Laboratory (<https://scgc.bigelow.org>). SAGs from phototrophic
126 (plastidic) and heterotrophic (aplastidic) cells were screened by PCR using universal
127 eukaryote DNA primers and taxonomically assigned. A total of 65 SAGs affiliated to 11
128 stramenopiles lineages (Table S1) were selected for sequencing. Sequence data is available at
129 ENA (<http://www.ebi.ac.uk/services/tara-oceans-data>) under the accession codes listed in
130 Table S2. Main sample-associated environmental data are reported in Table S3, and more
131 details can be found in PANGAEA (Pesant et al., 2015).

132 *SAG sequencing and assembly*

133 After purification of MDA products, 101 bp paired-end libraries were prepared from each
134 single cell as described in Alberti et al., 2017 and cells were independently sequenced on a
135 1/8th Illumina HiSeq lane at the Oregon Health & Science University (US) or at the National

136 Sequencing Center of Genoscope (France). Reads from SAGs were assembled using SPAdes
137 3.1 (Nurk et al., 2013). In all assemblies, contigs shorter than 500bp were discarded. Quality
138 profiles and basic statistics (genome size, number of contigs, N50, GC content) of each SAG
139 assemblies were generated with Quast (Gurevich, Saveliev, Vyahhi, & Tesler, 2013).
140 Estimations of genome recovery were done with BUSCO (Benchmarking Universal Single-
141 Copy Orthologs; Simão, Waterhouse, Ioannidis, Kriventseva, & Zdobnov, 2015).

142 *Detection and identification of viral signals in SAGs*

143 Putative viral sequences were retrieved from each assembled SAG using VirSorter v1.0.3
144 (Roux, Enault, Hurwitz, & Sullivan, 2015) with default parameters and both the *RefSeqABVir*
145 and *Virome* databases through the CyVerse Discovery Environment (Devisetty, Kennedy,
146 Sarando, Merchant, & Lyons, 2016). Contigs identified by VirSorter at all three levels of
147 confidence (from the more to the less confident predictions), categorized as viruses and
148 prophages, were used in subsequent analyses. Sequence similarity between identified full
149 length viral contigs was checked via pairwise BLASTn v2.2.28 (Altschul, Gish, Miller,
150 Myers, & Lipman, 1990). Contigs with sequence similarity >95%, coverage >80% and e-
151 value of $<10^{-5}$ were clustered together. Only one representative contig (i.e., the longest one)
152 for each non-redundant SAG-associated viral sequence of each cluster (here after unique
153 contig) was kept for further analysis.

154 Taxonomy of SAG-associated viral contigs was inferred using the webserver ViPTree
155 (Nishimura et al., 2017). Proteomic trees of each unique contig were generated based on
156 genome-wide sequence similarities computed by tBLASTx. A measure of genomic similarity
157 based on a normalized bit score of tBLASTx (S_G) was calculated against a set of reference
158 viral genomes database, the GenomeNet Virus-Host Database (Mihara et al., 2016). Since

159 MDA does not amplify RNA viruses, only ssDNA and dsDNA viruses and
160 virophage/satellites were considered in that analysis. SAG-associated viruses showing a S_G
161 >0.15 with a reference viral genome were assumed to belong to the same viral genus
162 (Nishimura et al., 2017).

163 Finally, protein-coding genes of each unique SAG-associated viral sequences were predicted
164 using Prodigal v2.6.3 (Hyatt et al., 2010) and annotated with BLASTp v2.7.1 (e-value 0.001,
165 max. 10 hits) using the NCBI's nr database (updated 09 Feb 2019).

166 *Biogeography of SAG-associated viral contigs assessed by fragment recruitment analysis.*

167 The global distribution of each unique SAG-associated virus was estimated by fragment
168 recruitment analysis against metagenomes from the Ocean Microbial Reference Gene Catalog
169 (OM-RGC; Sunagawa et al., 2015) were estimated using an approach similar to Swan et al.,
170 2013. A total of 128 metagenomes from two depths (surface and Deep Chlorophyll Maximum
171 [DCM]) targeting both $<0.22 \mu\text{m}$ ($n = 48$) and $0.22\text{-}3 \mu\text{m}$ ($n= 80$) size fractions were
172 analyzed. Metagenomic reads were prior randomly subsampled without replacement to the
173 minimum number of reads within each depth and size fraction using reformat.sh from bbtools
174 suite (<https://sourceforge.net/projects/bbmap/>). BLAST+ v2.7.1 was then used to recruit reads
175 from the OM-RGC database to each viral sequence ($n=64$) using default parameter values,
176 except for: -perc_identity 70 -evalue 0.0001. The percentage of unique recruits (~ 100 bp long
177 and $\geq 95\%$ identity) from each metagenome matching to each viral sequence was normalized
178 by viral sequence length. SAG-associated viral sequence abundances for each metagenome
179 were calculated from the BLAST output and plotted using custom R scripts.

180

181 *Identification of virophage contigs and detection/reconstruction of the mavirus integration*
182 *site*

183 The SAG-associated viral contig SV11, determined by ViPTree in the previous analyses, was
184 highly similar to the virophage genome Maverick-related virus (also referred as mavirus,
185 NC_015230, Fischer & Hackl, 2016), which share an evolutionary origin with a class of self-
186 synthesizing DNA transposons called Maverick/Polinton elements (Fischer & Suttle, 2011).
187 Mavirus was recently found integrated within the nuclear genome of the protist *Cafeteria*
188 *roenbergensis* in multiple sites, where the endogenous virophage genome (named
189 *Cafeteriavirus*-dependent mavirus and here referred as endogenous mavirus, KU052222) was
190 flanked on either side by terminal inverted repeats (TIRs) (Fischer & Hackl, 2016). However,
191 in comparison with the sequence of the endogenous mavirus, SV11 virophage sequence was
192 partially incomplete. To determine if the incomplete SAG-associated virophages were
193 potentially integrated within their respective host genome, we proceeded as follows. We first
194 identified putative virophage contigs that could have not been detected with VirSorter because
195 this automated tool was only applied on contigs >500bp, and requires a minimum of two
196 predicted genes per contig to identify it as viral. Consequently, for each SAG containing an
197 associated putative virophage sequence (within chrysophyte-G1 and MAST-3A), all contigs
198 (including fragments <500 bp) were searched by a BLASTn analysis against a manually
199 curated sequence of the endogenous mavirus including the TIRs sequences. For each SAG,
200 contigs with a minimum similarity of 95% and maximal e-value of 10^{-4} with the curated
201 endogenous mavirus genome were assumed to belong to the virophage genome and were
202 aligned to the primarily detected virophage contig using ClustalW (Larkin et al., 2007) as
203 implemented in the Geneious package v10.2.2 (Kearse et al., 2012). Then, to increase the
204 completeness of the SAG-associated virophage genome, a fragment recruitment analysis was

205 performed using BLASTn against all identified virophage contigs in each SAG and reads with
206 at least 99% identity and a maximal e-value of 10^{-4} were kept and assembled to the virophage
207 genome using the Geneious *de novo* assembler with a minimum overlap of 50bp and a
208 minimum identity of 95% (Kearse et al., 2012). Gene prediction of the obtained SAG-
209 associated virophage assemblies was done using Prodigal v2.6.3 and annotated with BLASTp
210 v2.2.28 (e-value 0.001, max. 10 hits) against NCBI's nr database (updated 06 Jun 2017).

211 *Phylogenetic and comparative genomic analysis of the virophages*

212 Phylogenetic analyses of the new SAG-associated virophages were performed with a set of
213 reference virophage sequences from the literature. These include the virophage sequences
214 isolated from cultures such as Sputnik (La Scola et al., 2008), Sputnik 2 (Desnues et al.,
215 2012), Sputnik 3 (Gaia et al., 2013), Zamilon (Gaia et al., 2014) and mavirus (Fischer &
216 Suttle, 2011), combined with sequences assembled from environmental surveys such as
217 Yellowstone Lake (YLSV1-4 (Zhou et al., 2013) and YLSV5-7 (Zhou et al., 2015), Qinghai
218 Lake (QLV, (Oh, Yoo, & Liu, 2016)), Dishui Lake (Dishui, (Gong et al., 2016)), Organic
219 Lake (OLV, (Yau et al., 2011)), Ace Lake (ALM, (Zhou et al., 2013)), Trout Bog epilimnion
220 and hypolimnion (TBE and TBH, Roux et al., 2017) and Mendota (Roux et al., 2017).

221 As proposed by Roux et al., 2017, phylogenetic trees were built based on a concatenated
222 alignment using four core genes (major capsid protein [MCP], minor capsid protein [mCP],
223 DNA packaging enzyme [ATPase], and cysteine protease [CysProt]) from all virophage
224 genomes, except for the virophage TBE_1002136, which lacked the ATPase. For this last,
225 only 3 genes were included in the multi-marker alignment. For each virophage core gene,
226 individual alignments were generated with MAFFT v7.305b (L-INS-I algorithm, (Katoh &
227 Standley, 2013)), automatically curated to remove all non-informative positions using trimAl

228 v1.2 (Capella-Gutierrez, Silla-Martinez, & Gabaldon, 2009) and evaluated for optimal amino
229 acid substitution models using ProtTest v3.4.2 (Darriba, Taboada, Doallo, & Posada, 2011).
230 The concatenation of the four core genes alignments was performed using a supermatrix
231 approach with a custom python script
232 (https://github.com/wrf/supermatrix/blob/master/add_taxa_to_align.py). Maximum-likelihood
233 trees of each four individual core genes alignments and the concatenated alignment were
234 constructed with RAxML v. 8.2.9 (Stamatakis, 2014) with 100 trees for both topology and
235 rapid bootstrap analyses, and using the evolutionary models LG+I+G+F (ATPase, CysProt
236 and mCP) and RtREV+I+G+F (MCP). Trees were generated using the ape (Paradis, Claude, &
237 Strimmer, 2004) and ggtree packages (Yu, Smith, Zhu, Guan, & Lam, 2017) in R 3.5.1. (R
238 Development Core Team, 2016), and rooted using QLV. To verify the topology of the trees,
239 bayesian phylogenies on each alignment were also generated with MrBayes v3.2.6 (2,000,000
240 generations; Ronquist et al., 2012).

241 Finally, whole-genome synteny comparisons between chrysophyte-G1 and MAST-3A SAG-
242 associated virophages and their closest published relatives (endogenous *Cafeteriavirus*-
243 dependent mavirus and Ace Lake mavirus) were performed with EasyFig v.2.2.2 (Sullivan,
244 Petty, & Beatson, 2011) using tBLASTx and filtering of small hits and annotations option.
245 Since all chrysophyte-G1 SAG-associated virophage are highly similar (mean identity of
246 99%), only one representative sequence (i.e., longest assembly) per stramenopile lineage are
247 displayed (AB233-L11 for chrysophyte-G1 and AA240-G22 for MAST-3A).

248

249 **Results**

250 *Detection of viral contigs in protist SAGs*

251 We used a total of 65 SAGs from photosynthetic and heterotrophic stramenopiles selected
252 from four *Tara* Oceans stations located in the Mediterranean Sea and Indian Ocean (Table
253 S3): 6 from two lineages of MAST-3 (clades A and F), 27 from three MAST-4 lineages
254 (clades A, C and E), 6 from a lineage of MAST-7 (clade A), 15 from three lineages of
255 Chrysophyceae (clades G1, H1 and H2), 4 from an uncultured clade of Dictyochophyceae and
256 7 affiliated to *Pelagomonas calceolata* (Table 1, Table S1). Using a relatively similar
257 sequencing depth (mean of 4.99 ± 0.81 Gbp), assembly sizes were variable among the SAGs,
258 averaging from 3.6 (± 2.8) Mbp in Dictyochophyceae to 11.0 (± 8.0) Mbp in MAST-3F
259 (considering contigs >500 bp; Table 1). The variation in assembly completeness was also
260 important, ranging on average from about 1% (in Pelagophyceae and Dictyochophyceae) to
261 10% (in MAST-3, MAST-4 and chrysophyte-G1; Table 1). Finally, the number of contigs
262 assembled and their respective N50 also varied among SAGs and stramenopile lineages
263 (Table 1).

264 We first investigated the presence of viral contigs in the 65 stramenopile SAGs assemblies,
265 identifying a total of 79 putative viral sequences in 37 SAGs (~57%) distributed among most
266 analyzed lineages, with the exception of *Pelagomonas calceolata* (Figure. 1a, Table S1). Only
267 two lineages (MAST-4C and MAST-4E) showed less than half of their cells harboring viral
268 contigs (Figure 1a, Table S1). Interestingly, a significant fraction of the SAG-associated
269 viruses (~50%) was found in cells affiliated to chrysophytes (8, 24 and 12 viral contigs in
270 chrysophyte-G1, -H1 and -H2, respectively; Figure 1a) with only one cell without any viral
271 sequence detected out of the 15 analyzed cells (Fig 1a). Furthermore, the isolated chrysophyte
272 cells were very rich in viral sequences, with up to 9 viral contigs retrieved in a single SAG of

273 chrysohyte-H1 (AA538_K19; Table S1). However, this was an exception, since in general
274 we detected from 1 to 3 viral contigs per cell (Table S1).

275 We next explored the uniqueness of the detected viral contigs based on a pairwise comparison
276 of their full-length sequences. Of the 79 viral sequences initially identified, we determined 64
277 non-redundant (i.e., unique) sequences (Table 2), ranging from 1 to 48.5 kbp in length
278 (median = 5.7 kbp; Table 2). From the 64, 61 were associated to a single stramenopile lineage
279 (~95%), and only 3 viral sequences (~5%) were either shared between two (SV11 and SV28)
280 or four lineages (SV2) (Figure 1b): SV2 was found in MAST-4 (clades A and E), MAST-7
281 and chrysohyte-H1, whereas SV11 was detected in chrysohyte-G1 and MAST-3A, and
282 SV28 in MAST-4 clades A and E (Table 2). With respect to the 61 viral contigs present in
283 only one specific lineage, about 98% of them were reported in only one specific cell (Figure
284 1b), with the exceptional case of SV51 present in triplicate in the same chrysohyte-H1 cell
285 (AA538_K19, Table S1). Only one viral contig (SV13) was found in two different
286 chrysohyte-H2 cells (Figure 1b, Table 2).

287 *Diversity and distribution of the SAG-associated viral sequences across the sunlit oceans*

288 The 64 unique viral sequences were compared with a set of reference viral genomes (Mihara
289 et al., 2016). On the basis of high genomic sequence similarity ($S_G > 0.15$), 7 of the 64 viruses
290 identified in the SAGs could be putatively assigned to four different viral families. These
291 viruses were two virophages (SV11, SV46), three viruses of *Phycodnaviridae* (SV27, SV35,
292 SV50), one virus of *Myoviridae* (SV48) and one virus of *Podoviridae* (SV64; Table 2). Other
293 viruses showed lower sequence similarities ($n = 48$; $S_G < 0.15$) or lacked detectable similarity
294 by tBLASTn ($n=9$) to reference viral genomes, thus being uncertain for their classification at
295 the genus level (Table 2). None of the assigned viral genomes were complete (or circular) but

296 one particular virus, SV11, which seemed nearly complete based on the similarity to a
297 reference genome. This virus of 15.5 kbp in length was highly similar (98.3%, $S_G = 0.96$) to
298 the Maverick-related virus genome (mavirus, GenBank accession number: NC_015230) and
299 likely belong to the virophage genus of *Mavirus* (*Lavidaviridae*; Table 2). The remaining
300 identified viral signals includes a set of short genome fragments (from 1 to 6.4 kbp) with
301 intermediate genomic similarities (40-53%) to either an unclassified virophage (SV46 with
302 YLV6), eukaryotic viruses (SV27 and SV50 with *Phycodnaviridae*), or phages (SV48 and
303 SV64) (Table 2). We further predicted protein-coding genes in the 64 unique viral sequences.
304 Of the total of 619 predicted genes (median = 6 predicted genes per SV; Table 2), about ~60%
305 (n= 363) had a close relative in the NCBI's nr database and 103 genes were related to
306 eukaryotic viral functions (Table S4).

307 In order to address the occurrence of these putative viruses in marine epipelagic waters, we
308 performed a fragment recruitment analysis of the viral signals in the *Tara* Oceans OM-RGC
309 database (Sunagawa et al., 2015). Our findings show that the viral contigs were found
310 preferentially at the DCM, and at the 0.2-3 μm size fraction rather than in the $<0.2 \mu\text{m}$ size
311 fraction (Figure 2 and Figure S2). Regarding their geographic distribution, the SAG-
312 associated viruses displayed some differences. On the one hand, some of them showed a
313 cosmopolitan distribution with different degrees of occurrence. For example, some viral
314 contigs (SV1 and SV2) show a high presence in all oceanic basins and in both size fractions,
315 while others (e.g., SV34) were highly present in the 0.2-3 μm size fraction but absent from the
316 $<0.2 \mu\text{m}$ size fraction (Figure 2). On the other hand, other SAG-associated viruses appeared to
317 be constrained to a lower number of oceanic basins, with some of them showing some
318 biogeography preferences (e.g., SV16 and SV32), whereas others were restricted to few
319 locations with a low presence (e.g., SV53, SV54 and SV55) (Figure 2 and Figure S2).

320 *Genome reconstruction and phylogenetic analysis of SAG-associated virophages*

321 We next focused on five SAG-associated viral contigs (the non-redundant SAG-associated
322 viral contig SV11), retrieved from one MAST-3A (AB240-G22) and from four different
323 chrysophyte-G1 cells (AB233-D06, AB233-L11, AB233-O05 and AB233-P23; Table 2 and
324 Table S1), which were highly similar to mavirus (i.e., Maverick-related virus; Table 2), an
325 endogenous virophage (“provirophage”) integrated in the genome of *Cafeteria roenbergensis*.
326 To the best of our knowledge, mavirus constitutes the only case of integration of a *Mavirus*
327 virophage in a protist genome revealed to date by a culture-based approach. However, the
328 virophage genomes identified in each SAG were incomplete compared to mavirus, noting the
329 remarkable lack of two genes coding for an integrase and an helicase, as well as the TIRs,
330 which indicate genome linearity and, therefore, a potential integration into the host genome
331 (Fischer & Hackl, 2016; Roux et al., 2017). After the identification of the putative virophage
332 contigs and a read recruitment analysis within each SAG, to increase the completeness of the
333 SAG-associated virophage genomes (see methods section), we were able to reconstruct the
334 entire SAG-associated virophage genomes of the five stramenopiles cells. This includes the
335 presence of both DNA replication genes and TIRs on either side of all SAG-associated
336 virophage genomes, confirming that SAG-associated mavirus genomes were linear and
337 potentially inserted in the stramenopile host genomes. For the reassembly, from 5 (AB233-
338 O05) to 14 contigs (AB240-G22), ranging from 0.2-0.3 to 7.2-15.5 kbp in length, were
339 necessary to reconstruct the 5 SAG-associated virophage genomes.

340 To better assess the phylogenetic position of these newly identified SAG-associated mavirus
341 genomes among the virophages, we established a concatenated marker tree using four
342 virophage core genes (mCP, MCP, ATPase and CysProt; Fig 3), including all the available

343 virophage genomes retrieved from culture, metagenomes and the five new SAG-associated
344 virophages. We found that the newly identified virophage sequences form a clade among the
345 genus *Mavirus* together with the mavirus virophage but distinct from the Ace Lake mavirus
346 (ALM (Zhou et al., 2013), a partial *Mavirus* genome retrieved from an environmental
347 sequencing survey) (Figure 3). Similar phylogenetic placements were found when each core
348 gene was analyzed separately (Fig. S1).

349 Finally, we compared the general genome organization of the identified SAG-associated
350 mavirus in MAST-3A (SV11_AB240_G22) and chrysophyte-G1 (SV11_AB233_L11) and
351 their closest published relatives, the endogenous mavirus and Ace Lake mavirus. As expected
352 from the previous analysis, the two SAG-associated mavirus displayed remarkable sequence
353 similarity with the endogenous mavirus integrated within the nuclear genome of *Cafeteria*
354 *roenbergensis* and exhibited clear differences with the Ace Lake mavirus (Figure 4). The
355 main differences between the two SAG-associated mavirus and the endogenous mavirus are
356 the presence of an extra gene coding for an unknown function (gene 11, 71 amino acids) in
357 the two SAG-associated mavirus and the absence in SV11_AB233_L11 mavirus of the gene
358 20 (152 amino acids, unknown function) of the endogenous mavirus genome (Figure 4).
359 Interestingly, we also retrieved an exon structure of one adjacent host gene of unknown
360 function (gene 22, 177 amino acids) in the MAST-3A genome (Figure 4). Although this
361 putative host sequence is relatively short (~ 1kbp), we were able to observe a significant
362 difference in its overall GC content compared with the mavirus sequence (60% vs 35%,
363 respectively; Figure 4).

364 **Discussion**

365 In this study, we used SCG to characterize potential virus-protist interactions. With the
366 exception of *Pelagomonas calceolata*, we found evidence of virus associations in almost all
367 studied protist cells. Indeed, the relatively high frequency of viral associations with protists
368 cells (~57%), retrieved from SAG assemblies with low genome recovery (<10%), suggests
369 that viral association levels are much higher. Same observations were previously made for
370 prokaryotic cells in marine environments (Labonté et al., 2015; Munson-McGee et al., 2018;
371 Roux et al., 2014), implying that (nearly) all microbial cells are susceptible to be infected or
372 to carry viruses. In the case of *Pelagomonas calceolata* cells, the lack of viral signals in their
373 respective SAGs is probably due to the very low assembly coverage ($0.5 \pm 0.4\%$) rather than
374 to the absence of any virus. We have not yet observed any significant correlation between
375 genome completeness and the number of viral contigs among SAGs (Table 1). Similar
376 findings were previously reported in bacterioplankton (Labonté et al., 2015), suggesting
377 that the probability to detect viruses among SAGs is independent of the retrieved host
378 genome assembly. The variation in SAG genome coverage may depend on intrinsic
379 properties of selected cells, their DNA integrity, as well as multiple displacement
380 amplification (MDA) biases (Pinard et al., 2006; Stepanauskas, 2012; Woyke et al., 2009).
381 Several methods have been developed to improve genome recovery of uncultured cells such
382 as using partial SAG assemblies to recruit metagenome reads and/or contigs (Saw et al.,
383 2015), sorting multiple natural cells to perform a targeted metagenomic analysis (Cuvelier et
384 al., 2010; Rinke et al., 2013; Vaultot et al., 2012) or co-assembling short reads from multiple
385 SAGs (Mangot et al., 2017; Seeleuthner et al., 2018). However, the application of these
386 approaches to characterize virus–host interactions will miss intraspecific genetic variability of
387 both actors. More recently, several new MDA-like methods, such as WGA-X (Stepanauskas
388 et al., 2017), TruePrime (Picher et al., 2016) or REPLI-g (Ahsanuddin et al., 2017) have been

389 developed for improving the genome recovery from single environmental cells (bacterial,
390 archaeal and protists) and viral particles with high GC-content genomes. Compared with
391 the conventional MDA, these amplification alternatives may provide a better genome
392 recovery of microbial taxa, including some not amenable to standard MDA (Stepanauskas et
393 al., 2017).

394 Using VirSorter (Roux et al., 2015) we were able to identify 64 unique viral contigs in 37
395 stramenopiles cells. We chose VirSorter over VirFinder (Ren, Ahlgren, Lu, Fuhrman, & Sun,
396 2017) because it has been shown that the later may misclassify eukaryotic sequences as viral.
397 Some other approaches have been recently developed to retrieve viral signals from (meta-)
398 genomic data, such as MARVEL (Amgarten, Braga, da Silva, & Setubal, 2018) and VirMiner
399 (Zheng et al., 2019), but they have been developed to detect viral genomes in prokaryotes.

400 From the 64 viral contigs retrieved in the protist cells, the narrow host range of these viruses
401 was remarkable given that >95% of the detected viral sequences (n=61) were specific to one
402 stramenopile lineage and just a few were shared between lineages (n=3). This is contrary to
403 previous findings on prokaryotic SAGs showing that nearly 50% of the detected viral types
404 were found in more than 2 lineages (Munson-McGee et al., 2018), suggesting that viruses
405 infecting protists are likely more specialist than viruses infecting prokaryotes. Furthermore,
406 while an important fraction (~54%) of cells with viral signals was associated to only one viral
407 sequence, we also retrieved several putative co-infections among the remaining cells, with up
408 to 7 unique (i.e., non-redundant) viral contigs in a single chrysophyte cell. Nonetheless, the
409 risk of a putative accidental co-sorting of a free viral particle with a protist cell during the
410 single-cell sorting process exists. To assess the risk of a possible “viral contamination”, we
411 estimated the frequency of such events based on previous estimates made on prokaryotic cells
412 (Labonté et al., 2015) by adapting the calculations to the cell size range of our studied

413 stramenopile cells (2-3 μm). We obtained that the frequency of free environmental viral
414 particles present in the cells' shade was less than 1 in 5,000. This reinforces the view that the
415 viruses detected in our study were truly and directly associated to the analyzed protist cells.
416 These associations may consist on i) lytic and/or temperate (i.e. non-lytic) viruses adsorbed in
417 the cell membrane, ii) a temperate virus or a virophage integrated into the host genome, iii) a
418 virus replicating inside the cell, iv) a grazed prokaryote or protist carrying a temperate virus
419 or with an active infection, or v) a predated free virus. A combination of these different
420 scenarios probably explains the high number of viral sequences detected in the chrysophyte
421 cells. For now, our current data set, including mostly fragments of viral sequences rather than
422 complete viral genomes, does not allow us to decipher which mode of virus-host association
423 prevail among the targeted protist cells.

424 Only 7 (~10%) viral contigs detected in protist cells were taxonomically assigned to known
425 viruses (Table 2), which include some close hits to viruses belonging to the *Phycodnaviridae*
426 family, known as a pathogen of marine eukaryotic algae (e.g., Brussaard, Short, Frederickson,
427 & Suttle, 2004; Derelle et al., 2008), and others to bacteriophages and cyanophages. This
428 suggests that a non-negligible part of the identified viral signals might come from putative
429 infected (bacterial and/or picoeukaryotic) preys grazed by the stramenopiles. Indeed, the
430 analyzed stramenopile lineages are mostly small free-living bacterivorous (Massana et al.,
431 2006; Piwosz, Wiktor, Niemi, Tatarek, & Michel, 2013), with some groups (e.g., MAST-4)
432 showing the ability to also eat picoalgae in grazing experiments (Massana et al., 2009).
433 Nevertheless, previous studies working with a subset of our SAGs (Mangot et al., 2017;
434 Seeleuthner et al., 2018) have shown that genes from bacteria and photosynthetic eukaryotes
435 only represent a very small fraction of the genome assemblies (< 0.3% of fragmented contigs
436 (Mangot et al., 2017)). A search for 16S rDNA genes in the SAGs where bacteriophages were

437 retrieved was unfruitful (data not shown), making difficult the association of these phages to
438 putative grazed bacteria. It is also possible that some of the detected viral signals come from
439 grazed viruses, since it is well-known that heterotrophic protists can graze on viruses
440 (Fuhrman, 1999; González & Suttle, 1993). However, another plausible explanation for the
441 identification of bacteriophages as closest hits to some SAGs associated virus is the
442 overrepresentation of bacteriophage genomes compared to viruses in reference databases
443 (Klingenberg, Aßhauer, Lingner, & Meinicke, 2013), which is supported by the low sequence
444 similarity between the viral signals and the bacteriophages sequences (Table 2). Although
445 taking all together it is difficult to elucidate which virus-host associations prevail among the
446 targeted protist cells, the geographic distribution of the viral signals supports the view that the
447 detected virus-protists associations reflect in many cases true interactions, because viral
448 signals coming from MAST-4A, MAST-4C and chrysophyte-H1 were ubiquitous (e.g. SV1,
449 SV7 and SV28), while those viral signals coming from MAST-4E, MAST-3A, MAST-3F and
450 chrysophyte-H2 were geographically constrained (e.g. SV12, SV32 and SV54) (Figure 2), in
451 agreement with the biogeography of these stramenopiles (Seeleuthner et al., 2018).

452 Some of the taxonomically assigned viral contigs were affiliated to known virophages and,
453 more particularly, in the case of SV11 to *Lavidaviridae* (Krupovic, Kuhn, & Fischer, 2016).
454 This virophage family, encompassing the two genera of *Mavirus* and *Sputnikvirus*, comprises
455 obligate parasites of giant DNA viruses of the *Mimiviridae* family (Fischer & Hackl, 2016).
456 Furthermore, virophages encode integrase genes, and proviropages have been reported in the
457 nuclear genome of the marine alga *Bigeloviella natans* (Blanc, Gallot-Lavallée, & Maumus,
458 2015), and the protozoan *Cafeteria roenbergensis* (Fischer & Hackl, 2016). Proviropages
459 putatively act as a host defense mechanism against giant viruses, in which some cells are
460 sacrificed to protect their kin (Blanc et al., 2015; Fischer & Suttle, 2011). In this study, we

461 identified the presence of endogenous mavirus virophages in the assembly of five cells
462 affiliated to chrysohyte-G1 and MAST-3A. These SAG-mavirus are highly similar to the
463 *Cafeteriavirus*-dependent mavirus, a parasite of the giant *Cafeteria roenbergensis* virus
464 (CroV) (Fischer, Allen, Wilson, & Suttle, 2010) integrated within the genome of *Cafeteria*
465 *roenbergensis* (Fischer & Hackl, 2016). The presence of TIRs in the SV11_AB233_L11 and
466 SV11_AB240_G22 virophages, as well as the exon structure of a putative adjacent host gene
467 in the SV11_AB240_G22 sequence (Figure 4), suggests the putative integration of the
468 mavirus in the host genome. This is also confirmed by the lack of any CroV signal in our
469 assembly, whose presence is incompatible with a virophage in its lysogenic stage (Fischer &
470 Hackl, 2016). This finding constitutes the first report of the presence of a putative
471 provirophage isolated from environmental samples using SCG. Only slight differences were
472 observed between the different provirophage genomes, located notably at genomic regions of
473 low conservation (gene 11 in the two SAG-associated mavirus). Little is known about the
474 importance of mavirus provirophage in protist populations as its study is limited to few cases.
475 It is somewhat surprising that the same mavirus virophage was found in three
476 phylogenetically distant lineages (chrysohyte-G1, MAST-3A and *C. roenbergensis*),
477 pointing to a global and important ecological role of virophages in protist populations.
478 Mavirus host cell recognition is carried out through specific receptor interactions, while
479 Sputnik is done through phagocytosis of a composite of the virophage and the giant virus they
480 parasitize (Duponchel & Fischer, 2019). Therefore, a possible explanation for finding mavirus
481 in the different lineages, is that the capsid proteins are evolutionary conserved and have
482 evolved independently of the giant virus infecting the host cell. On the contrary, although the
483 host cells from Sputnik and Zamilon are phylogenetically closer, these virophages may have
484 co-evolved with their corresponding giant virus. This hypothesis is supported by the finding

485 that Sputnik can infect the groups A, B and C of the Mimiviridae group while Zamilon is
486 unable to infect the group A (mimi- and mamavirus) (Gaia et al., 2014). Virophages are
487 repeatedly detected in genomic studies, with different gene content and abundance profiles,
488 likely suggesting that they occupy different ecological niches (Desnues & Raoult, 2012; Roux
489 et al., 2017; Yau et al., 2011). Although the role of virophages in protist populations is still
490 enigmatic, they may play a role in regulating the giant virus population dynamics and virus-
491 host interactions, influencing the ecosystem function and probably the whole microbial food
492 web in aquatic environments (Desnues & Raoult, 2012). Our findings provide new insights
493 into the potential importance of mavirus in the ecology of marine protists, and reinforce the
494 need for more studies to elucidate the role of these fascinating viruses in the environment.

495

496 In summary, this work shows the benefits of single-cell genomics to increase our
497 understanding of virus-host associations in natural protist communities. Although our
498 knowledge of the marine viral diversity is constantly expanding since the development of
499 metagenomics (Coutinho et al., 2017; Mizuno et al., 2013; Paez-Espino et al., 2016), it has
500 been estimated that the majority (63-93%) of viral sequences in marine metagenomes are not
501 represented in public databases (Hurwitz & Sullivan, 2013), emphasizing the need for further
502 isolation, characterization and sequencing of specific marine viruses (Middelboe & Brussaard,
503 2017). A minute fraction of protist viruses is annotated to date (~100 sequenced genomes
504 (~0.6% of all viral genomes) in NCBI Genome database (July 2018), explaining the majority
505 of unassigned viral sequences in our study. Thus, in addition to the ever-increasing
506 knowledge on viral diversity by metagenomic approaches, the incorporation of SAG analysis
507 will allow the specific matching of viruses and their hosts as well as to determine the host
508 range of individual viruses without cultivation. Our findings suggest that protist cells are

509 susceptible to interact with predominantly specialist viruses and hint to the potential
510 importance of provirophages in protist populations.

511 **Acknowledgments**

512 This work was supported by the Spanish projects MEFISTO (CTM2013-43767-P, MINECO),
513 ALLFLAGS (CTM2016-75083-R, MINECO) and INTERACTOMICS (CTM2015-69936-P,
514 MINECO), and the EU project SINGEK (H2020-MSCA-ITN-2015-675752). YMC was
515 supported by a FPI Spanish fellowship (BES-2014-067849). JFM was beneficiary of a Marie
516 Curie Fellowship (PIEF-GA-2012-331190, EU). LFB was beneficiary of a Marie Curie
517 Fellowship (H2020-MSCA-ITN-2015-675752, EU). RL was supported by a Ramón y Cajal
518 fellowship (RYC-2013-12554, MINECO, Spain). HO was supported by JSPS/KAKENHI
519 (No. 18H02279), and Scientific Research on Innovative Areas from the Ministry of
520 Education, Culture, Science, Sports and Technology (MEXT) of Japan (Nos. 16H06429,
521 16K21723, 16H06437). OJ was supported by The French Government ‘Investissement
522 d’Avenir’ programs Oceanomics (ANR-11-BTBR-0008) and FRANCE GENOMIQUE
523 (ANR-10-INBS-09). MS was supported by a Viera y Clavijo contract funded by the ACIISI
524 and the ULPGC. Computing resources were obtained through the MARBITS platform at the
525 ICM-CSIC. We are grateful to Michael E. Sieracki, Patrick Wincker and Colomban de
526 Vargas, members of the *Tara* Oceans consortium, who initiated and designed the sampling
527 and sequencing experiments of protist SAGs. We thank Nigel Grimsley, Matthias Fischer and
528 Simon Roux for their help on the early stage of the analysis of the SAG-associated viroplage
529 sequences. We are also grateful to Pablo Sánchez for his advices on computing the fragment
530 recruitment analysis. We finally thank the *Tara* Oceans consortium, people, and sponsors who

531 supported the *Tara* Oceans Expedition (<http://www.embl.de/tara-oceans/>) for making the data
532 accessible. This is the contribution number XXX of the *Tara* Oceans Expedition 2009-2013.

533

534 **References**

- 535 Ahsanuddin, S., Afshinnakoo, E., Gandara, J., Hakyemezoglu, M., Bezdán, D., Minot, S., ...
536 Mason, C. E. (2017). Assessment of REPLI-g multiple displacement whole genome
537 amplification (WGA) techniques for metagenomic applications. *Journal of Biomolecular*
538 *Techniques*, 28(1), 46–55. doi:10.7171/jbt.17-2801-008
- 539 Alberti, A., Poulain, J., Engelen, S., Labadie, K., Romac, S., Ferrera, I., ... Wincker, P.
540 (2017). Viral to metazoan marine plankton nucleotide sequences from the *Tara* Oceans
541 expedition. *Scientific Data*, 4, 170093. doi:10.1038/sdata.2017.93
- 542 Allers, E., Moraru, C., Duhaime, M. B., Beneze, E., Solonenko, N., Barrero-Canosa, J., ...
543 Sullivan, M. B. (2013). Single-cell and population level viral infection dynamics
544 revealed by phageFISH, a method to visualize intracellular and free viruses.
545 *Environmental Microbiology*, 15(8), 2306–2318. doi:10.1111/1462-2920.12100
- 546 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local
547 alignment search tool. *Journal of Molecular Biology*, 215(3), 403–410.
548 doi:10.1016/S0022-2836(05)80360-2
- 549 Amgarten, D., Braga, L. P. P., da Silva, A. M., & Setubal, J. C. (2018). MARVEL, a tool for
550 prediction of bacteriophage sequences in metagenomic bins. *Frontiers in Genetics*, 9,
551 304. doi:10.3389/fgene.2018.00304
- 552 Anderson, R. E., Brazelton, W. J., & Baross, J. A. (2011). Using CRISPRs as a metagenomic
553 tool to identify microbial hosts of a diffuse flow hydrothermal vent viral assemblage.
554 *FEMS Microbiology Ecology*, 77(1), 120–133. doi:10.1111/j.1574-6941.2011.01090.x
- 555 Baran, N., Goldin, S., Maidanik, I., & Lindell, D. (2018). Quantification of diverse virus
556 populations in the environment using the polony method. *Nature Microbiology*, 3(1), 62–
557 72. doi:10.1038/s41564-017-0045-y
- 558 Berg Miller, M. E., Yeoman, C. J., Chia, N., Tringe, S. G., Angly, F. E., Edwards, R. A., ...
559 White, B. A. (2012). Phage-bacteria relationships and CRISPR elements revealed by a
560 metagenomic survey of the rumen microbiome. *Environmental Microbiology*, 14(1),
561 207–227. doi:10.1111/j.1462-2920.2011.02593.x
- 562 Bhattacharya, D., Price, D. C., Yoon, H. S., Yang, E. C., Poulton, N. J., Andersen, R. A., &
563 Das, S. P. (2012). Single cell genome analysis supports a link between phagotrophy and
564 primary plastid endosymbiosis. *Scientific Reports*, 2(1), 356. doi:10.1038/srep00356
- 565 Blanc, G., Gallot-Lavallée, L., & Maumus, F. (2015). Provirophages in the *Bigeloviella*
566 genome bear testimony to past encounters with giant viruses. *Proceedings of the*
567 *National Academy of Sciences*, 112(38), E5318–E5326. doi:10.1073/pnas.1506469112
- 568 Bolduc, B., Wirth, J. F., Mazurie, A., & Young, M. J. (2015). Viral assemblage composition
569 in Yellowstone acidic hot springs assessed by network analysis. *The ISME Journal*,
570 9(10), 2162–2177. doi:10.1038/ismej.2015.28
- 571 Breitbart, M. (2012). Marine viruses: truth or dare. *Annual Review of Marine Science*, 4(1),
572 425–448. doi:10.1146/annurev-marine-120709-142805

573 Breitbart, M., Bonnain, C., Malki, K., & Sawaya, N. A. (2018). Phage puppet masters of the
574 marine microbial realm. *Nature Microbiology*, *3*(7), 754–766. doi:10.1038/s41564-018-
575 0166-y

576 Brum, J. R., Ignacio-Espinoza, J. C., Roux, S., Doulier, G., Acinas, S. G., Alberti, A., ...
577 Sullivan, M. B. (2015). Patterns and ecological drivers of ocean viral communities.
578 *Science*, *348*(6237), 1261498–1261498. doi:10.1126/science.1261498

579 Brum, J. R., & Sullivan, M. B. (2015). Rising to the challenge: accelerated pace of discovery
580 transforms marine virology. *Nature Reviews Microbiology*, *13*(3), 147–159.
581 doi:10.1038/nrmicro3404

582 Brussaard, C P D, Short, S. M., Frederickson, C. M., & Suttle, C. A. (2004). Isolation and
583 phylogenetic analysis of novel viruses infecting the phytoplankton *Phaeocystis globosa*
584 (Prymnesiophyceae). *Applied and Environmental Microbiology*, *70*(6), 3700–3705.
585 doi:10.1128/AEM.70.6.3700-3705.2004

586 Brussaard, Corina P D, Wilhelm, S. W., Thingstad, F., Weinbauer, M. G., Bratbak, G.,
587 Haldal, M., ... Wommack, K. E. (2008). Global-scale processes with a nanoscale drive:
588 the role of marine viruses. *The ISME Journal*, *2*(6), 575–578. doi:10.1038/ismej.2008.31

589 Brüssow, H., & Hendrix, R. W. (2002). Phage genomics: small is beautiful. *Cell*, *108*(1), 13–
590 16. doi:10.1016/S0092-8674(01)00637-7

591 Capella-Gutierrez, S., Silla-Martinez, J. M., & Gabaldon, T. (2009). trimAl: a tool for
592 automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*,
593 *25*(15), 1972–1973. doi:10.1093/bioinformatics/btp348

594 Chow, C. E. T., Winget, D. M., White, R. A., Hallam, S. J., & Suttle, C. A. (2015).
595 Combining genomic sequencing methods to explore viral diversity and reveal potential
596 virus-host interactions. *Frontiers in Microbiology*, *6*, 1–15.
597 doi:10.3389/fmicb.2015.00265

598 Coutinho, F. H., Silveira, C. B., Gregoracci, G. B., Thompson, C. C., Edwards, R. A.,
599 Brussaard, C. P. D., ... Thompson, F. L. (2017). Marine viruses discovered via
600 metagenomics shed light on viral strategies throughout the oceans. *Nature*
601 *Communications*, *8*, 15955. doi:10.1038/ncomms15955

602 Cuvelier, M. L., Allen, A. E., Monier, A., McCrow, J. P., Messie, M., Tringe, S. G., ...
603 Worden, A. Z. (2010). Targeted metagenomics and ecology of globally important
604 uncultured eukaryotic phytoplankton. *Proceedings of the National Academy of Sciences*,
605 *107*(33), 14679–14684. doi:10.1073/pnas.1001665107

606 Danovaro, R., Corinaldesi, C., Dell’Anno, A., Fuhrman, J. A., Middelburg, J. J., Noble, R. T.,
607 & Suttle, C. A. (2011). Marine viruses and global climate change. *FEMS Microbiology*
608 *Reviews*, *35*(6), 993–1034. doi:10.1111/j.1574-6976.2010.00258.x

609 Darriba, D., Taboada, G. L., Doallo, R., & Posada, D. (2011). ProtTest 3: fast selection of
610 best-fit models of protein evolution. *Bioinformatics*, *27*(8), 1164–1165.
611 doi:10.1093/bioinformatics/btr088

612 del Campo, J., & Massana, R. (2011). Emerging diversity within Chrysophytes,
613 Choanoflagellates and Bicosoecids based on molecular surveys. *Protist*, *162*(3), 435–
614 448. doi:10.1016/j.protis.2010.10.003

615 Deng, L., Gregory, A., Yilmaz, S., Poulos, B. T., Hugenholtz, P., & Sullivan, M. B. (2012).
616 Contrasting life strategies of viruses that infect photo- and heterotrophic bacteria, as
617 revealed by viral tagging. *MBio*, *3*(6), e00373-12. doi:10.1128/mBio.00373-12

618 Deng, L., Ignacio-Espinoza, J. C., Gregory, A. C., Poulos, B. T., Weitz, J. S., Hugenholtz, P.,
619 & Sullivan, M. B. (2014). Viral tagging reveals discrete populations in *Synechococcus*
620 viral genome sequence space. *Nature*, *513*(7517), 242–245. doi:10.1038/nature13459

621 Derelle, E., Ferraz, C., Escande, M.-L., Eychenié, S., Cooke, R., Piganeau, G., ... Grimsley,
622 N. (2008). Life-cycle and genome of OtV5, a large DNA virus of the pelagic marine
623 unicellular green alga *Ostreococcus tauri*. *PLoS ONE*, 3(5), e2250.
624 doi:10.1371/journal.pone.0002250

625 Desnues, C., La Scola, B., Yutin, N., Fournous, G., Robert, C., Azza, S., ... Raoult, D.
626 (2012). Provirophages and transpovirons as the diverse mobilome of giant viruses.
627 *Proceedings of the National Academy of Sciences*, 109(44), 18078–18083.
628 doi:10.1073/pnas.1208835109

629 Desnues, C., & Raoult, D. (2012). Virophages question the existence of satellites. *Nature*
630 *Reviews Microbiology*, 10(3), 234–234. doi:10.1038/nrmicro2676-c3

631 Devisetty, U. K., Kennedy, K., Sarando, P., Merchant, N., & Lyons, E. (2016). Bringing your
632 tools to CyVerse Discovery Environment using Docker. *F1000Research*, 5, 1442.
633 doi:10.12688/f1000research.8935.1

634 Duponchel, S., & Fischer, M. G. (2019). Viva lavidaviruses! Five features of virophages that
635 parasitize giant DNA viruses. *PLOS Pathogens*, 15(3), e1007592.
636 doi:10.1371/journal.ppat.1007592

637 Fischer, M. G., Allen, M. J., Wilson, W. H., & Suttle, C. A. (2010). Giant virus with a
638 remarkable complement of genes infects marine zooplankton. *Proceedings of the*
639 *National Academy of Sciences*, 107(45), 19508–19513. doi:10.1073/pnas.1007615107

640 Fischer, M. G., & Hackl, T. (2016). Host genome integration and giant virus-induced
641 reactivation of the virophage mavirus. *Nature*, 540(7632), 288–291.
642 doi:10.1038/nature20593

643 Fischer, M. G., & Suttle, C. A. (2011). A virophage at the origin of large DNA transposons.
644 *Science*, 332(6026), 231–234. doi:10.1126/science.1199412

645 Fuhrman, J. A. (1999). Marine viruses and their biogeochemical and ecological effects.
646 *Nature*, 399(6736), 541–548. doi:10.1038/21119

647 Gaia, M., Benamar, S., Boughalmi, M., Pagnier, I., Croce, O., Colson, P., ... La Scola, B.
648 (2014). Zamilon, a novel virophage with Mimiviridae host specificity. *PLoS ONE*, 9(4),
649 e94923. doi:10.1371/journal.pone.0094923

650 Gaia, M., Pagnier, I., Campocasso, A., Fournous, G., Raoult, D., & La Scola, B. (2013).
651 Broad spectrum of Mimiviridae virophage allows its isolation using a Mimivirus
652 reporter. *PLoS ONE*, 8(4), e61912. doi:10.1371/journal.pone.0061912

653 Gong, C., Zhang, W., Zhou, X., Wang, H., Sun, G., Xiao, J., ... Wang, Y. (2016). Novel
654 virophages discovered in a freshwater lake in China. *Frontiers in Microbiology*, 7, 5.
655 doi:10.3389/fmicb.2016.00005

656 González, J., & Suttle, C. (1993). Grazing by marine nanoflagellates on viruses and virus-
657 sized particles: ingestion and digestion. *Marine Ecology Progress Series*, 94, 1–10.
658 doi:10.3354/meps094001

659 Gurevich, A., Saveliev, V., Vyahhi, N., & Tesler, G. (2013). QUASt: quality assessment tool
660 for genome assemblies. *Bioinformatics*, 29(8), 1072–1075.
661 doi:10.1093/bioinformatics/btt086

662 Heywood, J. L., Sieracki, M. E., Bellows, W., Poulton, N. J., & Stepanauskas, R. (2011).
663 Capturing diversity of marine heterotrophic protists: one cell at a time. *The ISME*
664 *Journal*, 5(4), 674–684. doi:10.1038/ismej.2010.155

665 Hurwitz, B. L., & Sullivan, M. B. (2013). The Pacific Ocean Virome (POV): a marine viral
666 metagenomic dataset and associated protein clusters for quantitative viral ecology. *PLoS*
667 *ONE*, 8(2), e57355. doi:10.1371/journal.pone.0057355

668 Hyatt, D., Chen, G.-L., LoCascio, P. F., Land, M. L., Larimer, F. W., & Hauser, L. J. (2010).

669 Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC*
670 *Bioinformatics*, 11(1), 119. doi:10.1186/1471-2105-11-119

671 Jover, L. F., Effler, T. C., Buchan, A., Wilhelm, S. W., & Weitz, J. S. (2014). The elemental
672 composition of virus particles: implications for marine biogeochemical cycles. *Nature*
673 *Reviews Microbiology*, 12(7), 519–528. doi:10.1038/nrmicro3289

674 Karsenti, E., Acinas, S. G., Bork, P., Bowler, C., De Vargas, C., Raes, J., ... Wincker, P.
675 (2011). A holistic approach to marine eco-systems biology. *PLoS Biology*, 9(10),
676 e1001177. doi:10.1371/journal.pbio.1001177

677 Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version
678 7: improvements in performance and usability. *Molecular Biology and Evolution*, 30(4),
679 772–780. doi:10.1093/molbev/mst010

680 Kearse, M., Moir, R., Wilson, A., Stones-Havas, S., Cheung, M., Sturrock, S., ... Drummond,
681 A. (2012). Geneious Basic: an integrated and extendable desktop software platform for
682 the organization and analysis of sequence data. *Bioinformatics*, 28(12), 1647–1649.
683 doi:10.1093/bioinformatics/bts199

684 Klingenberg, H., Aßhauer, K. P., Lingner, T., & Meinicke, P. (2013). Protein signature-based
685 estimation of metagenomic abundances including all domains of life and viruses.
686 *Bioinformatics*, 29(8), 973–980. doi:10.1093/bioinformatics/btt077

687 Krabberød, A., Bjorbækmo, M., Shalchian-Tabrizi, K., & Logares, R. (2017). Exploring the
688 oceanic microeukaryotic interactome with metaomics approaches. *Aquatic Microbial*
689 *Ecology*, 79(1), 1–12. doi:10.3354/ame01811

690 Krupovic, M., Kuhn, J. H., & Fischer, M. G. (2016). A classification system for virophages
691 and satellite viruses. *Archives of Virology*, 161(1), 233–247. doi:10.1007/s00705-015-
692 2622-9

693 La Scola, B., Desnues, C., Pagnier, I., Robert, C., Barrassi, L., Fournous, G., ... Raoult, D.
694 (2008). The virophage as a unique parasite of the giant mimivirus. *Nature*, 455(7209),
695 100–104. doi:10.1038/nature07218

696 Labonté, J. M., Swan, B. K., Poulos, B., Luo, H., Koren, S., Hallam, S. J., ... Stepanauskas,
697 R. (2015). Single-cell genomics-based analysis of virus–host interactions in marine
698 surface bacterioplankton. *The ISME Journal*, 9(11), 2386–2399.
699 doi:10.1038/ismej.2015.48

700 Larkin, M. A., Blackshields, G., Brown, N. P., Chenna, R., McGettigan, P. A., McWilliam,
701 H., ... Higgins, D. G. (2007). Clustal W and Clustal X version 2.0. *Bioinformatics*,
702 23(21), 2947–2948. doi:10.1093/bioinformatics/btm404

703 Lin, Y.-C. C., Campbell, T., Chung, C.-C. C., Gong, G.-C. C., Chiang, K.-P. P., & Worden,
704 A. Z. (2012). Distribution patterns and phylogeny of marine stramenopiles in the North
705 Pacific Ocean. *Applied and Environmental Microbiology*, 78(9), 3387–3399.
706 doi:10.1128/AEM.06952-11

707 Logares, R., Audic, S., Santini, S., Pernice, M. C., de Vargas, C., & Massana, R. (2012).
708 Diversity patterns and activity of uncultured marine heterotrophic flagellates unveiled
709 with pyrosequencing. *The ISME Journal*, 6(10), 1823–1833. doi:10.1038/ismej.2012.36

710 Mangot, J.-F., Logares, R., Sánchez, P., Latorre, F., Seeleuthner, Y., Mondy, S., ... Massana,
711 R. (2017). Accessing the genomic information of unculturable oceanic picoeukaryotes by
712 combining multiple single cells. *Scientific Reports*, 7, 41498. doi:10.1038/srep41498

713 Massana, R. (2011). Eukaryotic picoplankton in surface oceans. *Annual Review of*
714 *Microbiology*, 65(1), 91–110. doi:10.1146/annurev-micro-090110-102903

715 Massana, R., del Campo, J., Sieracki, M. E., Audic, S., & Logares, R. (2014). Exploring the
716 uncultured microeukaryote majority in the oceans: reevaluation of ribogroups within

717 stramenopiles. *The ISME Journal*, 8(4), 854–866. doi:10.1038/ismej.2013.204

718 Massana, R., Terrado, R., Forn, I., Lovejoy, C., & Pedrós-Alió, C. (2006). Distribution and
719 abundance of uncultured heterotrophic flagellates in the world oceans. *Environmental*
720 *Microbiology*, 8(9), 1515–1522. doi:10.1111/j.1462-2920.2006.01042.x

721 Massana, R., Unrein, F., Rodríguez-Martínez, R., Forn, I., Lefort, T., Pinhassi, J., & Not, F.
722 (2009). Grazing rates and functional diversity of uncultured heterotrophic flagellates.
723 *The ISME Journal*, 3(5), 588–596. doi:10.1038/ismej.2008.130

724 Middelboe, M., & Brussaard, C. (2017). Marine viruses: key players in marine ecosystems.
725 *Viruses*, 9(10), 302. doi:10.3390/v9100302

726 Mihara, T., Nishimura, Y., Shimizu, Y., Nishiyama, H., Yoshikawa, G., Uehara, H., ... Ogata,
727 H. (2016). Linking virus genomes with host taxonomy. *Viruses*, 8(3), 66.
728 doi:10.3390/v8030066

729 Mizuno, C. M., Rodriguez-Valera, F., Kimes, N. E., & Ghai, R. (2013). Expanding the marine
730 virosphere using metagenomics. *PLoS Genetics*, 9(12), e1003987.
731 doi:10.1371/journal.pgen.1003987

732 Munn, C. B. (2006). Viruses as pathogens of marine organisms—from bacteria to whales.
733 *Journal of the Marine Biological Association of the UK*, 86(03), 453–467.
734 doi:10.1017/S002531540601335X

735 Munson-McGee, J. H., Peng, S., Dewerff, S., Stepanauskas, R., Whitaker, R. J., Weitz, J. S.,
736 & Young, M. J. (2018). A virus or more in (nearly) every cell: ubiquitous networks of
737 virus–host interactions in extreme environments. *The ISME Journal*, 12(7), 1706–1714.
738 doi:10.1038/s41396-018-0071-7

739 Nishimura, Y., Watai, H., Honda, T., Mihara, T., Omae, K., Roux, S., ... Yoshida, T. (2017).
740 Environmental viral genomes shed new light on virus-host interactions in the ocean.
741 *MSphere*, 2(2). doi:10.1128/mSphere.00359-16

742 Nurk, S., Bankevich, A., Antipov, D., Gurevich, A. A., Korobeynikov, A., Lapidus, A., ...
743 Pevzner, P. A. (2013). Assembling single-cell genomes and mini-metagenomes from
744 chimeric MDA products. *Journal of Computational Biology*, 20(10), 714–737.
745 doi:10.1089/cmb.2013.0084

746 Oh, S., Yoo, D., & Liu, W.-T. (2016). Metagenomics reveals a novel virophage population in
747 a tibetan mountain lake. *Microbes and Environments*, 31(2), 173–177.
748 doi:10.1264/jsme2.ME16003

749 Paez-Espino, D., Eloie-Fadrosch, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann, M.,
750 Mikhailova, N., ... Kyrpides, N. C. (2016). Uncovering Earth’s virome. *Nature*,
751 536(7617), 425–430. doi:10.1038/nature19094

752 Paradis, E., Claude, J., & Strimmer, K. (2004). APE: Analyses of Phylogenetics and
753 Evolution in R language. *Bioinformatics*, 20(2), 289–290.
754 doi:10.1093/bioinformatics/btg412

755 Pesant, S., Not, F., Picheral, M., Kandels-Lewis, S., Le Bescot, N., Gorsky, G., ... Searson, S.
756 (2015). Open science resources for the discovery and analysis of *Tara* Oceans data.
757 *Scientific Data*, 2, 150023. doi:10.1038/sdata.2015.23

758 Picher, Á. J., Budeus, B., Wafzig, O., Krüger, C., García-Gómez, S., Martínez-Jiménez, M. I.,
759 ... Schneider, A. (2016). TruePrime is a novel method for whole-genome amplification
760 from single cells based on TthPrimPol. *Nature Communications*, 7(1), 13296.
761 doi:10.1038/ncomms13296

762 Pinard, R., de Winter, A., Sarkis, G. J., Gerstein, M. B., Tartaro, K. R., Plant, R. N., ...
763 Leamon, J. H. (2006). Assessment of whole genome amplification-induced bias through
764 high-throughput, massively parallel whole genome sequencing. *BMC Genomics*, 7(1),

765 216. doi:10.1186/1471-2164-7-216

766 Piwosz, K., Wiktor, J. M., Niemi, A., Tatarek, A., & Michel, C. (2013). Mesoscale
767 distribution and functional diversity of picoeukaryotes in the first-year sea ice of the
768 Canadian Arctic. *The ISME Journal*, 7(8), 1461–1471. doi:10.1038/ismej.2013.39

769 R Development Core Team. (2016). R: a language and environment for statistical computing.
770 Retrieved from <https://www.r-project.org/>

771 Rappé, M. S., & Giovannoni, S. J. (2003). The uncultured microbial majority. *Annual Review*
772 *of Microbiology*, 57(1), 369–394. doi:10.1146/annurev.micro.57.030502.090759

773 Ren, J., Ahlgren, N. A., Lu, Y. Y., Fuhrman, J. A., & Sun, F. (2017). VirFinder: a novel k-
774 mer based tool for identifying viral sequences from assembled metagenomic data.
775 *Microbiome*, 5(1), 69. doi:10.1186/s40168-017-0283-5

776 Rinke, C., Schwientek, P., Sczyrba, A., Ivanova, N. N., Anderson, I. J., Cheng, J.-F., ...
777 Woyke, T. (2013). Insights into the phylogeny and coding potential of microbial dark
778 matter. *Nature*, 499(7459), 431–437. doi:10.1038/nature12352

779 Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D. L., Darling, A., Höhna, S., ...
780 Huelsenbeck, J. P. (2012). MrBayes 3.2: efficient Bayesian phylogenetic inference and
781 model choice across a large model space. *Systematic Biology*, 61(3), 539–542.
782 doi:10.1093/sysbio/sys029

783 Roux, S., Brum, J. R., Dutilh, B. E., Sunagawa, S., Duhaime, M. B., Loy, A., ... Sullivan, M.
784 B. (2016). Ecogenomics and potential biogeochemical impacts of globally abundant
785 ocean viruses. *Nature*, 537(7622), 689–693. doi:10.1038/nature19366

786 Roux, S., Chan, L.-K., Egan, R., Malmstrom, R. R., McMahon, K. D., & Sullivan, M. B.
787 (2017). Ecogenomics of virophages and their giant virus hosts assessed through time
788 series metagenomics. *Nature Communications*, 8(1), 858. doi:10.1038/s41467-017-
789 01086-2

790 Roux, S., Enault, F., Hurwitz, B. L., & Sullivan, M. B. (2015). VirSorter: mining viral signal
791 from microbial genomic data. *PeerJ*, 3, e985. doi:10.7717/peerj.985

792 Roux, S., Hawley, A. K., Torres Beltran, M., Scofield, M., Schwientek, P., Stepanauskas, R.,
793 ... Sullivan, M. B. (2014). Ecology and evolution of viruses infecting uncultivated
794 SUP05 bacteria as revealed by single-cell- and meta-genomics. *ELife*, 3, e03125.
795 doi:10.7554/eLife.03125

796 Roy, R. S., Price, D. C., Schliep, A., Cai, G., Korobeynikov, A., Yoon, H. S., ...
797 Bhattacharya, D. (2015). Single cell genome analysis of an uncultured heterotrophic
798 stramenopile. *Scientific Reports*, 4(1), 4780. doi:10.1038/srep04780

799 Saw, J. H., Spang, A., Zaremba-Niedzwiedzka, K., Juzokaite, L., Dodsworth, J. A.,
800 Murugapiran, S. K., ... Ettema, T. J. G. (2015). Exploring microbial dark matter to
801 resolve the deep archaeal ancestry of eukaryotes. *Philosophical Transactions of the*
802 *Royal Society of London. Series B, Biological Sciences*, 370(1678), 20140328.
803 doi:10.1098/rstb.2014.0328

804 Seeleuthner, Y., Mondy, S., Lombard, V., Carradec, Q., Pelletier, E., Wessner, M., ...
805 Wincker, P. (2018). Single-cell genomics of multiple uncultured stramenopiles reveals
806 underestimated functional diversity across oceans. *Nature Communications*, 9(1), 310.
807 doi:10.1038/s41467-017-02235-3

808 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E. M. (2015).
809 BUSCO: assessing genome assembly and annotation completeness with single-copy
810 orthologs. *Bioinformatics*, 31(19), 3210–3212. doi:10.1093/bioinformatics/btv351

811 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis
812 of large phylogenies. *Bioinformatics*, 30(9), 1312–1313.

813 doi:10.1093/bioinformatics/btu033

814 Stepanauskas, R. (2012). Single cell genomics: an individual look at microbes. *Current*

815 *Opinion in Microbiology*, 15(5), 613–620. doi:10.1016/j.mib.2012.09.001

816 Stepanauskas, R., Fergusson, E. A., Brown, J., Poulton, N. J., Tupper, B., Labonté, J. M., ...

817 Lubys, A. (2017). Improved genome recovery and integrated cell-size analyses of

818 individual uncultured microbial cells and viral particles. *Nature Communications*, 8(1),

819 84. doi:10.1038/s41467-017-00128-z

820 Sullivan, M. J., Petty, N. K., & Beatson, S. A. (2011). Easyfig: a genome comparison

821 visualizer. *Bioinformatics*, 27(7), 1009–1010. doi:10.1093/bioinformatics/btr039

822 Sunagawa, S., Coelho, L. P., Chaffron, S., Kultima, J. R., Labadie, K., Salazar, G., ...

823 Velayoudon, D. (2015). Structure and function of the global ocean microbiome. *Science*,

824 348(6237), 1261359–1261359. doi:10.1126/science.1261359

825 Suttle, C. A. (2005). Viruses in the sea. *Nature*, 437(7057), 356–361.

826 doi:10.1038/nature04160

827 Swan, B. K., Tupper, B., Sczyrba, A., Lauro, F. M., Martinez-Garcia, M., Gonzalez, J. M., ...

828 Stepanauskas, R. (2013). Prevalent genome streamlining and latitudinal divergence of

829 planktonic bacteria in the surface ocean. *Proceedings of the National Academy of*

830 *Sciences*, 110(28), 11463–11468. doi:10.1073/pnas.1304246110

831 Tadmor, A. D., Ottesen, E. A., Leadbetter, J. R., & Phillips, R. (2011). Probing individual

832 environmental bacteria for viruses by using microfluidic digital PCR. *Science*,

833 333(6038), 58–62. doi:10.1126/science.1200758

834 Troell, K., Hallström, B., Divne, A.-M., Alsmark, C., Arrighi, R., Huss, M., ... Bertilsson, S.

835 (2016). *Cryptosporidium* as a testbed for single cell genome characterization of

836 unicellular eukaryotes. *BMC Genomics*, 17(1), 471. doi:10.1186/s12864-016-2815-y

837 Vannier, T., Leconte, J., Seeleuthner, Y., Mondy, S., Pelletier, E., Aury, J.-M., ... Jaillon, O.

838 (2016). Survey of the green picoalga *Bathycoccus* genomes in the global ocean.

839 *Scientific Reports*, 6(1), 37900. doi:10.1038/srep37900

840 Vaultot, D., Lepère, C., Toulza, E., De la Iglesia, R., Poulain, J., Gaboyer, F., ... Piganeau, G.

841 (2012). Metagenomes of the picoalga *Bathycoccus* from the Chile coastal upwelling.

842 *PLoS ONE*, 7(6), e39648. doi:10.1371/journal.pone.0039648

843 Weitz, J., & Wilhelm, S. (2012). Ocean viruses and their effects on microbial communities

844 and biogeochemical cycles. *FI000 Biology Reports*, 4, 17. doi:10.3410/B4-17

845 Woyke, T., Xie, G., Copeland, A., González, J. M., Han, C., Kiss, H., ... Stepanauskas, R.

846 (2009). Assembling the marine metagenome, one cell at a time. *PLoS ONE*, 4(4), e5299.

847 doi:10.1371/journal.pone.0005299

848 Yau, S., Lauro, F. M., DeMaere, M. Z., Brown, M. V., Thomas, T., Raftery, M. J., ...

849 Cavicchioli, R. (2011). Virophage control of antarctic algal host-virus dynamics.

850 *Proceedings of the National Academy of Sciences*, 108(15), 6163–6168.

851 doi:10.1073/pnas.1018221108

852 Yoon, H. S., Price, D. C., Stepanauskas, R., Rajah, V. D., Sieracki, M. E., Wilson, W. H., ...

853 Bhattacharya, D. (2011). Single-cell genomics reveals organismal interactions in

854 uncultivated marine protists. *Science*, 332(6030), 714–717. doi:10.1126/science.1203163

855 Yu, G., Smith, D. K., Zhu, H., Guan, Y., & Lam, T. T.-Y. (2017). ggtree: an R package for

856 visualization and annotation of phylogenetic trees with their covariates and other

857 associated data. *Methods in Ecology and Evolution*, 8(1), 28–36. doi:10.1111/2041-

858 210X.12628

859 Zheng, T., Li, J., Ni, Y., Kang, K., Misiakou, M.-A., Imamovic, L., ... Panagiotou, G. (2019).

860 Mining, analyzing, and integrating viral signals from metagenomic data. *Microbiome*,

861 7(1), 42. doi:10.1186/s40168-019-0657-y
862 Zhou, J., Sun, D., Childers, A., McDermott, T. R., Wang, Y., & Liles, M. R. (2015). Three
863 novel virophage genomes discovered from Yellowstone Lake metagenomes. *Journal of*
864 *Virology*, 89(2), 1278–1285. doi:10.1128/JVI.03039-14
865 Zhou, J., Zhang, W., Yan, S., Xiao, J., Zhang, Y., Li, B., ... Wang, Y. (2013). Diversity of
866 virophages in metagenomic data sets. *Journal of Virology*, 87(8), 4225–4236.
867 doi:10.1128/JVI.03398-12
868

869 **Data accessibility**

870 SAG sequence data are available at ENA under the accession codes listed in Table S2.

871 **Author contribution**

872 DV and RM conceived the study. OJ sequenced the SAGs, RL performed the SAG
873 assemblies, and YMC and JFM performed the SAG bioinformatic and data analysis. MS
874 contributed to the analytic design and LFB, HO and MK contributed with additional data
875 analyses. YMC, JFM, MS, DV and RM interpreted the results, and YMC and JFM wrote the
876 manuscript with inputs from all co-authors. YMC and JFM should be considered joint first
877 author.

878 **ORCID**

879 *Yaiza M. Castillo* <https://orcid.org/0000-0002-1319-2975>

880 **Supporting information**

881 Additional supporting information may be found in the online version of this article.

882 **FIGURE LEGEND**

883 **FIGURE 1** Occurrence and specificity of viral contigs in 65 marine stramenopiles SAGs. (a)
884 Barplots show the number of SAGs with or without any viral contig detected in their
885 assembly. For each lineage, the total number of SAG-associated viral contigs retrieved in

886 SAGs are indicated on top of each bar. **(b)** Pie charts display the percentage of viral contigs
887 present in only one or shared among 2 or 4 lineages (upper left corner), and the percentage of
888 SAGs that shared a viral contig for those that were lineage-specific (lower right corner).
889 Chryso-G1, Chryso-H1, Chryso-H2, Dictyo and Pelago correspond to the chrysophyte clades
890 G1, H1 and H2, Dictyochophyceae and *Pelagomonas calceolata*, respectively.

891 **FIGURE 2** Biogeographical distribution of SAG-associated viruses, as determined by
892 metagenomic fragment recruitment. Viral contigs are shown in the y-axis and epipelagic
893 metagenome stations along the x-axis. The scale bar indicates the percentage of read
894 sequences recruited normalized of aligned metagenome sequences with alignments ≥ 100 bp
895 long and $\geq 95\%$ identity, normalized by the length of each SAG-associated virus sequence.
896 Results for metagenomes from the $<0.2 \mu\text{m}$ (left panels) and $0.2\text{-}3 \mu\text{m}$ size fractions (right
897 panels) were displayed for both surface (upper panels) and DCM stations (lower panels).
898 Stations where metagenomes were not available are shown in grey (No Data). Color bars
899 represent the different oceanic basins, abbreviations: North Pacific Ocean (NP), South Pacific
900 Ocean (SP), North Atlantic Ocean (NA), South Atlantic Ocean (SA), Southern Ocean (SO),
901 Mediterranean Sea (M), Red Sea (RS) and Indian Ocean (IND).

902 **FIGURE 3** Phylogenetic placement of the new putative SAG-associated mavirus among
903 virophages. The tree topology was inferred from a maximum-likelihood analysis of a
904 concatenated alignment of four core genes (minor [mCP] and major [MCP] capsids proteins,
905 DNA packaging enzyme [ATPase] and Cysteine Protease [CysProt]). Bayesian posterior
906 probabilities (BPP) and bootstrap percentages (BS) are provided at each node (BPP/BS) when
907 support values were higher than 0.7 and 70%, respectively. Black dots indicate maximal
908 support for both posterior probabilities (1.0) and maximum-likelihood bootstraps (100%) at
909 the respective nodes. The five new SAG-associated virophages are highlighted in bold. The
910 origin (culture, metagenome sequencing or SAG) and genome type (linear with TIRs, circular
911 or partial) of each virophage genome are pointed out in the tree. Abbreviated names for
912 virophages are detailed in the Materials and Methods section.

913 **FIGURE 4** Comparison of the SAG-associated mavirus genomes and their closest known
914 relatives. Linear genomic maps show synteny between the mavirus genomes found in
915 chrysophyte-G1 and MAST-3A SAGs (SV11_AB233_L11 and SV11_AB240_G22,
916 respectively) and their closest published relatives, endogenous mavirus and ALM (Zhou et al.,
917 2013). When present, TIRs and exon structures of putative adjacent host genes are displayed
918 to highlight the putative integration of mavirus genomes within their respective host genomes.
919 The main differences between the two SAG-associated mavirus and the endogenous mavirus
920 are indicated with asterisks (\dagger : presence of an extra coding gene [gene 11], \ddagger : absence of
921 coding gene [gene 20]). Additionally, a GC content plot based on a 100 bp sliding window is
922 shown for SV11_AB240_G22.

TABLE 1 General characteristics (mean (\pm standard error)) of the 65 draft stramenopiles SAGs obtained by single-cell genomics

Group	Name	Number of cells	Sequencing depth (Gbp)	Assembly size (Mbp)	Total number of contigs	BUSCO completeness (%)	GC content (%)	N50 (kbp)
Chrysophyceae	Chrysophyte-G1	4	5.5 (\pm 0.5)	9.3 (\pm 5.2)	3,597 (\pm 2,009)	11.6 (\pm 8.2)	40.2 (\pm 0.2)	5.1 (\pm 0.4)
	Chrysophyte-H1 [†]	8	4 (\pm 0.7)	4.0 (\pm 2.1)	1,425 (\pm 518)	6.1 (\pm 3.4)	45.1 (\pm 0.8)	8.6 (\pm 3.4)
	Chrysophyte-H2 [†]	3	4.2 (\pm 1.0)	4.3 (\pm 2.4)	1,928 (\pm 1,073)	3.3 (\pm 1.9)	47.7 (\pm 1.7)	4.1 (\pm 0.2)
Dictyochophyceae	unc. dictyochophyte	4	4.6 (\pm 0.1)	3.6 (\pm 2.8)	15,567 (\pm 948)	1.0 (\pm 1.0)	46.8 (\pm 2.5)	4.2 (\pm 1.4)
MAST-3	MAST-3A	4	5.1 (\pm 0.5)	7.5 (\pm 1.9)	2,272 (\pm 409)	11.4 (\pm 3.4)	42.5 (\pm 0.3)	8.6 (\pm 1.1)
	MAST-3F	2	5.4 (\pm 1.1)	11.0 (\pm 8.0)	3,576 (\pm 2,414)	11.7 (\pm 9.8)	34.1 (\pm 0.3)	7.7 (\pm 0.2)
MAST-4	MAST-4A	14	5 (\pm 1.8)	10.2 (\pm 4.9)	3,195 (\pm 1,393)	12.6 (\pm 7.5)	33.0 (\pm 1.0)	9.3 (\pm 2.7)
	MAST-4C	4	5.4 (\pm 0.6)	8.3 (\pm 2.3)	2,389 (\pm 579)	11.8 (\pm 3.8)	40.3 (\pm 0.2)	14.0 (\pm 1.3)
	MAST-4E	9	4.7 (\pm 0.8)	6.7 (\pm 2.6)	1,928 (\pm 607)	8.6 (\pm 3.9)	44.0 (\pm 0.7)	9.3 (\pm 1.9)
MAST-7	MAST-7A	6	5.5 (\pm 1.3)	5.6 (\pm 3.2)	2,002 (\pm 1,233)	3.8 (\pm 2.1)	44.7 (\pm 4.8)	7.0 (\pm 3.2)
Pelagophyceae	<i>Pelagomonas calceolata</i>	7	5.6 (\pm 0.5)	8.1 (\pm 0.8)	271 (\pm 186)	0.5 (\pm 1.0)	47.5 (\pm 7.6)	13.0 \pm

Abbreviations: SAG, Single Amplified Genome; MAST, Marine Stramenopiles; unc., uncultured; BUSCO, Benchmarking Universal Single-Copy Orthologs; N50, length of the shortest contig from the minimal set of contig representing 50% of the assembly size.

[†] The 18S rRNA genes of these chrysophytes-H SAGs clustered into two distinct lineages (clades -H1 and -H2).

TABLE 2 Summary and taxonomic assignment of the 64 SAG-associated viral contigs

SAG-associated viral contig [†]				Taxonomic assignment on the GenomeNet Virus–Host Database (Mihara et al., 2016)				
Viral contig	SAG lineage (number of SAGs)	Sequence length (kbp)	Number of genes	Best viral group hit (GenBank accession number)	Viral family	Known host group	SG [‡]	Similarity (%)
SV1	MAST-4A (1)	48.5	44	<i>Cellulophaga</i> phage (KC821612)	Podoviridae	Bacteroidetes	0.06	40.8
SV2	MAST-4A (1), MAST-4E (1), MAST-7 (4), Chryso-H1 (1)	22.8	48	<i>Prochlorococcus</i> phage (NC_006883)	Myoviridae	Cyanobacteria	0.07	43.2
SV3	MAST-4A (1)	22.1	25	YSLV5 (NC_028269)	Unclassified viroplage	N/D	< 0.01	31.9
SV4	Chryso-H2 (1)	21.2	25	<i>Synechococcus</i> phage (NC_026928)	Myoviridae	Cyanobacteria	< 0.01	42.0
SV5	Chryso-H1 (1)	20.5	20	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.04	42.4
SV6	MAST-3A (1)	18.9	19	<i>Paramecium bursaria</i> <i>Chlorella</i> virus (NC_009898)	Phycodnaviridae	Ciliophora	0.04	39.2
SV7	Chryso-H1 (1)	17.9	32	<i>Pseudomonas</i> phage (NC_028980)	Siphoviridae	Gammaproteobacteria	0.06	58.7
SV8	MAST-4E (1)	16.7	19	<i>Phaeocystis globosa</i> virus viroplage (NC_021333)	Unclassified viroplage	Haptophyta	0.01	42.3
SV9	Chryso-H1 (1)	16.7	22	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.07	41.5
SV10	Chryso-H1 (1)	16.2	19	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.08	40.0
SV11	Chryso-G1 (4), MAST-3A (1)	15.5	18	Maverick-related virus (mavirus, NC_015230)	Lavidaviridae	Bicosoecophyceae	0.96	98.3
SV12	MAST-4C (1)	15.0	14	<i>Rhodothermus</i> phage (NC_004735)	Myoviridae	Bacteroidetes	0.02	39.6
SV13	Chryso-H2 (2)	14.3	11	<i>Chrysochromulina ericina</i> virus (NC_028094)	Phycodnaviridae	Haptophyta	0.04	66.6
SV14	MAST-4A (1)	13.1	14	-	-	-	-	-
SV15	MAST-4A (1)	12.7	11	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.02	37.6
SV16	MAST-3A (1)	12.6	18	Yellowstone lake phycodnavirus (NC_028110)	Phycodnaviridae	N/D	0.09	52.6
SV17	Chryso-G1 (1)	10.2	4	<i>Mycobacterium</i> phage (NC_028662)	Podoviridae	Actinobacteria	0.03	38.5
SV18	Chryso-H1 (1)	10.0	10	<i>Bacillus</i> phage (NC_006945)	Tectiviridae	Firmicutes	< 0.01	36.1
SV19	MAST-4E (1)	9.9	5	<i>Vibrio</i> phage (NC_021529)	Myoviridae	Gammaproteobacteria	0.06	49.0
SV20	MAST-7 (1)	9.8	8	<i>Cronobacter</i> phage (NC_019398)	Myoviridae	Gammaproteobacteria	0.02	51.2
SV21	Chryso-H1 (1)	8.6	8	<i>Phaeocystis globosa</i> virus viroplage (NC_021333)	Unclassified viroplage	Haptophyta	0.02	35.3
SV22	Chryso-G1 (1)	7.5	11	<i>Rhodothermus</i> phage (NC_015286)	Myoviridae	Cyanobacteria	0.05	49.7
SV23	MAST-7 (1)	7.4	7	<i>Acanthocystis turfacea</i> <i>Chlorella</i> virus (NC_008724)	Phycodnaviridae	Chlorophyta	0.08	54.9
SV24	Chryso-H1 (1)	7.2	5	-	-	-	-	-
SV25	Dictyo (1)	6.8	8	<i>Pseudomonas</i> phage (NC_026600)	Myoviridae	Gammaproteobacteria	0.01	43.1
SV26	Chryso-H1 (1)	6.7	9	-	-	-	-	-
SV27	Chryso-H2 (1)	6.4	8	Aureococcus anophagefferens virus (NC_024697)	Phycodnaviridae	Pelagophyceae	1.0	52.4
SV28	MAST-4A (2), MAST-4E (2)	5.9	6	<i>Cellulophaga</i> phage (KC821612)	Podoviridae	Bacteroidetes	0.07	44.2
SV29	Chryso-H1 (1)	5.8	7	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.14	41.0
SV30	Dictyo (1)	5.8	8	<i>Ostreococcus tauri</i> virus (NC_010191)	Phycodnaviridae	Chlorophyta	0.02	42.9
SV31	Chryso-H1 (1)	5.8	10	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.08	40.9
SV32	MAST-3A (1)	5.7	6	<i>Phaeocystis globosa</i> virus (NC_021312)	Phycodnaviridae	Haptophyta	0.07	45.4
SV33	Chryso-G1 (1)	5.6	3	<i>Anomala cuprea</i> entomovirus (NC_023426)	Poxviridae	Arthropoda	0.1	41.5
SV34	MAST-4A (1)	5.6	2	<i>Aureococcus anophagefferens</i> virus (NC_024697)	Phycodnaviridae	Pelagophyceae	< 0.01	39.1
SV35	MAST-3F (1)	5.4	7	Yellowstone lake phycodnavirus (NC_028110)	Phycodnaviridae	N/D	0.3	52.9
SV36	Chryso-H1 (1)	5.2	5	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.07	42.5
SV37	Chryso-H1 (1)	5.1	10	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.08	39.8
SV38	Chryso-H2 (1)	4.7	7	<i>Enterobacteria</i> phage (NC_005066)	Myoviridae	Gammaproteobacteria	0.02	39.7
SV39	Chryso-H1 (1)	4.6	3	<i>Phaeocystis globosa</i> virus viroplage (NC_021333)	Unclassified viroplage	Haptophyta	0.04	35.0
SV40	Dictyo (1)	4.6	4	<i>Enterobacteria</i> phage (NC_019526)	Myoviridae	Gammaproteobacteria	0.13	44.6
SV41	Chryso-H2 (1)	4.4	7	YSLV5 (NC_028269)	Unclassified viroplage	N/D	0.04	41.3
SV42	MAST-3A (1)	4.3	6	<i>Campylobacter</i> phage (NC_027997)	Myoviridae	Epsilonproteobacteria	0.02	30.3
SV43	Chryso-H1 (1)	4.1	6	YSLV7 (NC_028257)	Unclassified viroplage	N/D	0.04	41.1
SV44	MAST-3A (1)	4.1	5	<i>Aureococcus anophagefferens</i> virus (NC_024697)	Phycodnaviridae	Pelagophyceae	0.03	32.8
SV45	MAST-4A (1)	3.7	3	<i>Erwinia</i> phage (HQ728263)	Myoviridae	Gammaproteobacteria	0.02	40.4
SV46	Chryso-H1 (1)	3.7	3	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.17	40.3
SV47	Chryso-H1 (1)	3.1	5	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.09	42.3
SV48	Chryso-G1 (1)	3.0	2	Escherichia phage (NC_025447)	Myoviridae	Gammaproteobacteria	0.2	41.4
SV49	Chryso-H2 (1)	3.0	4	-	-	-	-	-
SV50	MAST-7 (1)	2.9	3	Ectocarpus siliculosus virus (NC_002687)	Phycodnaviridae	Phaeophyceae	0.2	46.1
SV51	Chryso-H1 (1)	2.9	5	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.1	46.5
SV52	Chryso-H1 (1)	2.9	4	-	-	-	-	-
SV53	Chryso-H2 (1)	2.8	4	-	-	-	-	-
SV54	Chryso-H2 (1)	2.8	4	-	-	-	-	-
SV55	Chryso-H2 (1)	2.7	4	-	-	-	-	-
SV56	Chryso-H1 (1)	2.5	4	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.04	42.3
SV57	MAST-4A (1)	2.5	4	<i>Synechococcus</i> phage (NC_015289)	Myoviridae	Cyanobacteria	0.03	31.1
SV58	MAST-7 (1)	2.3	3	<i>Enterobacteria</i> phage (NC_012740)	Myoviridae	Gammaproteobacteria	0.05	44.9
SV59	Chryso-H1 (1)	2.3	4	YSLV7 (NC_028257)	Unclassified viroplage	N/D	0.09	43.3
SV60	Chryso-H1 (1)	2.1	4	YSLV6 (NC_028270)	Unclassified viroplage	N/D	0.06	43.8
SV61	MAST-4A (1)	2.0	4	-	-	-	-	-
SV62	Chryso-H2 (1)	2.0	4	<i>Microcystis</i> phage (NC_029002)	Myoviridae	Cyanobacteria	0.04	40.4
SV63	Chryso-H1 (1)	1.8	4	YSLV5 (NC_028269)	Unclassified viroplage	N/D	0.1	47.0
SV64	MAST-4A (1)	1.0	3	Planktothrix phage (NC_016564)	Podoviridae	Cyanobacteria	0.2	49.1

Abbreviations: SAG, Single Amplified Genome; MAST, Marine Stramenopiles; Chryso, Chrysophyte; Dictyo, Dictyochophyceae; YSLV, Yellowstone Lake viroplage.

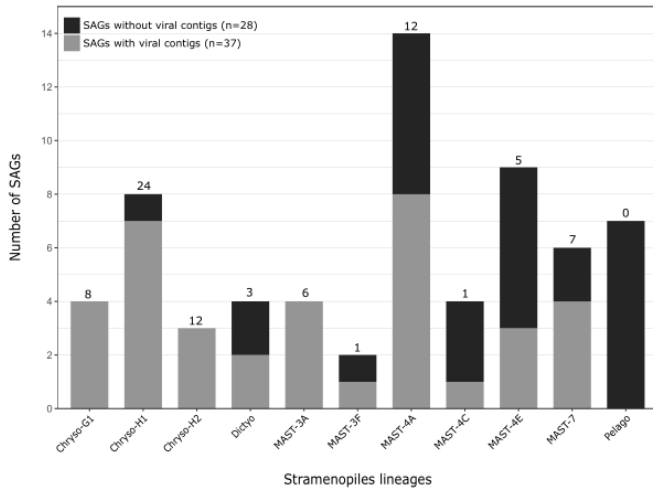
[†] Statistics are computed on the longest sequence when SAG-associated viral sequences were retrieved in several cell.

[‡] SG tBLASTx score. In bold, SAG-associated viral sequence that can be affiliated to the same genus level than their reference best hit (SG > 0.15).

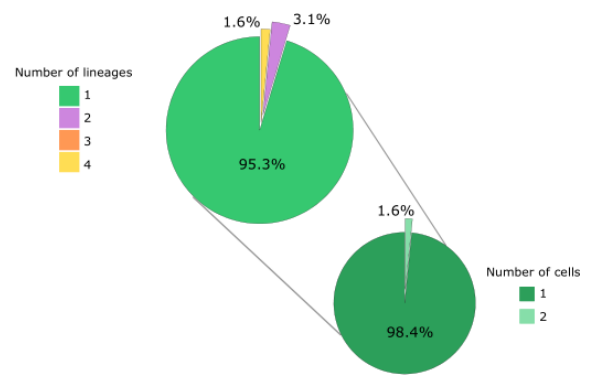
N/D Non Determined. Sequence were isolated from environmental surveys.

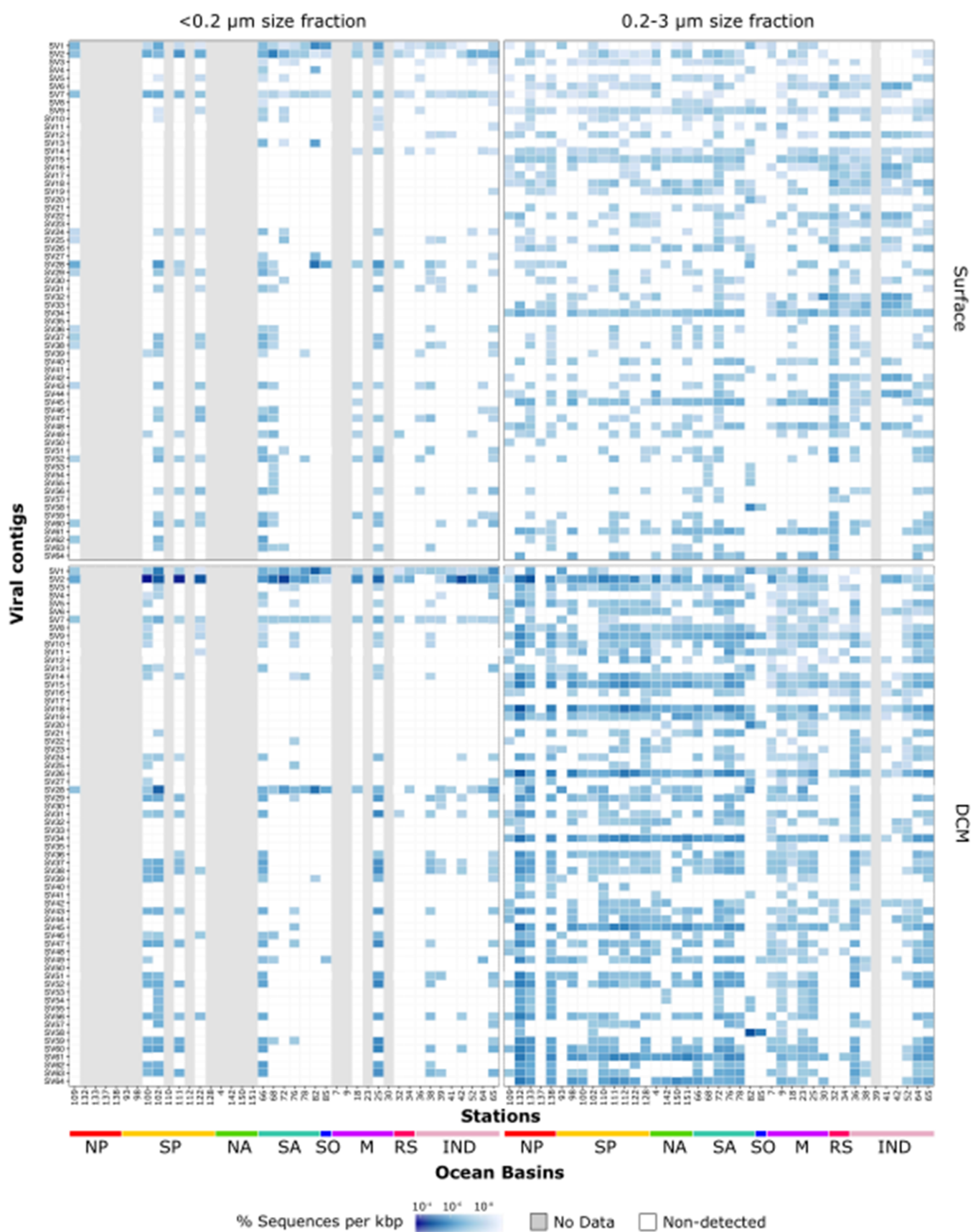
Viral signal sequence without taxonomic assignment are shown by the symbol (-)

(a)



(b)





Virophage origin Genome type

