

1 **Eukaryotic virus composition can predict the efficiency of carbon export in the**  
2 **global ocean**

3 Hiroto Kaneko<sup>1,+</sup>, Romain Blanc-Mathieu<sup>1,2,+</sup>, Hisashi Endo<sup>1</sup>, Samuel Chaffron<sup>3,4</sup>,  
4 Tom O. Delmont<sup>4,5</sup>, Morgan Gaia<sup>4,5</sup>, Nicolas Henry<sup>6</sup>, Rodrigo Hernández-Velázquez<sup>1</sup>,  
5 Canh Hao Nguyen<sup>1</sup>, Hiroshi Mamitsuka<sup>1</sup>, Patrick Forterre<sup>7</sup>, Olivier Jaillon<sup>4,5</sup>,  
6 Colomban de Vargas<sup>6</sup>, Matthew B. Sullivan<sup>8</sup>, Curtis A. Suttle<sup>9</sup>, Lionel Guidi<sup>10</sup> and  
7 Hiroyuki Ogata<sup>1,\*</sup>

8 + Equal contribution

9 \* Corresponding author

10 **Affiliations:**

11 1: Bioinformatics Center, Institute for Chemical Research, Kyoto University,  
12 Gokasho, Uji, Kyoto 611-0011, Japan

13 2: Laboratoire de Physiologie Cellulaire & Végétale, CEA, Univ. Grenoble Alpes,  
14 CNRS, INRA, IRIG, Grenoble, France.

15 3: Université de Nantes, CNRS UMR 6004, LS2N, F-44000 Nantes, France.

16 4: Research Federation (FR2022) *Tara* Oceans GO-SEE, Paris, France

17 5: Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, 91000  
18 Evry, France

19 6: Sorbonne Universités, CNRS, Laboratoire Adaptation et Diversité en Milieu Marin,  
20 Station Biologique de Roscoff, 29680 Roscoff, France

21 7: Institut Pasteur, Department of Microbiology, 25 rue du Docteur Roux, 75015,  
22 Paris, France

23 8: Department of Microbiology and Department of Civil, Environmental and Geodetic  
24 Engineering, Ohio State University, Columbus, OH, United States of America  
25 9: Departments of Earth, Ocean & Atmospheric Sciences, Microbiology &  
26 Immunology, and Botany, and the Institute for the Oceans and Fisheries, University  
27 of British Columbia, Vancouver, BC, V6T 1Z4, Canada  
28 10: Sorbonne Université, CNRS, Laboratoire d’Océanographie de Villefranche, LOV,  
29 F-06230 Villefranche-sur-mer, France  
30

## 31 **Summary**

32 The biological carbon pump, in which carbon fixed by photosynthesis is exported to  
33 the deep ocean through sinking, is a major process in Earth's carbon cycle. The  
34 proportion of primary production that is exported is termed the carbon export  
35 efficiency (CEE). Based on in-lab or regional scale observations, viruses were  
36 previously suggested to affect the CEE (i.e., viral "shunt" and "shuttle"). In this study,  
37 we tested associations between viral community composition and CEE measured at a  
38 global scale. A regression model based on relative abundance of viral marker genes  
39 explained 67% of the variation in CEE. Viruses with high importance in the model  
40 were predicted to infect ecologically important hosts. These results are consistent with  
41 the view that the viral shunt and shuttle functions at a large scale and further imply  
42 that viruses likely act in this process in a way dependent on their hosts and ecosystem  
43 dynamics.

## 44 **Introduction**

45 A major process in the global cycling of carbon is the oceanic biological carbon pump  
46 (BCP), an organism-driven process by which atmospheric carbon (*i.e.*, CO<sub>2</sub>) is  
47 transferred and sequestered to the ocean interior and seafloor for periods ranging from  
48 centuries to hundreds of millions of years. Between 15% and 20% of net primary  
49 production (NPP) is exported out of the euphotic zone, with 0.3% of fixed carbon  
50 reaching the seafloor annually (Zhang et al., 2018). However, there is wide variation  
51 in estimates of the proportion of primary production in the surface ocean that is  
52 exported to depth, ranging from 1% in the tropical Pacific to 35-45% during the North  
53 Atlantic bloom (Buesseler and Boyd, 2009). As outlined below, many factors affect  
54 the BCP.

55 Of planktonic organisms living in the upper layer of the ocean, diatoms  
56 (Tréguer et al., 2018) and zooplankton (Turner, 2015) have been identified as  
57 important contributors to the BCP in nutrient-replete oceanic regions. In the  
58 oligotrophic ocean, cyanobacteria, collodarians (Lomas and Moran, 2011), diatoms  
59 (Agusti et al., 2015; Karl et al., 2012; Leblanc et al., 2018), and other small (pico- to  
60 nano-) plankton (Lomas and Moran, 2011) have been implicated in the BCP.

61 Sediment trap studies suggest that ballasted aggregates of plankton with biogenic  
62 minerals contribute to carbon export to the deep sea (Iversen and Ploug, 2010; Klaas  
63 and Archer, 2002). The BCP comprises three processes: carbon fixation, export, and  
64 remineralization. As these processes are governed by complex interactions between  
65 numerous members of planktonic communities (Zhang et al., 2018), the BCP is  
66 expected to involve various organisms, including viruses (Zimmerman et al., 2019a).

67 Viruses have been suggested to regulate the efficiency of the BCP. Lysis of  
68 host cells by viruses releases cellular material in the form of dissolved organic matter

69 (DOM), which fuels the microbial loop and enhances respiration and secondary  
70 production (Gobler et al., 1997; Weitz et al., 2015). This process, coined “viral shunt  
71 (Wilhelm and Suttle, 1999)”, can reduce the carbon export efficiency (CEE) because  
72 it increases the retention of nutrients and carbon in the euphotic zone and prevents  
73 their transfer to higher trophic levels as well as their export from the euphotic zone to  
74 the deep sea (Fuhrman, 1999; Weitz et al., 2015). However, an alternative process is  
75 also considered, in which viruses contribute to the vertical carbon export (Weinbauer,  
76 2004). For instance, a theoretical study proposed that the CEE increases if viral lysis  
77 augments the ratio of exported carbon relative to the primary production-limiting  
78 nutrients (nitrogen and phosphorous) (Suttle, 2007). Laboratory experimental studies  
79 reported that cells infected with viruses form larger particles (Peduzzi and Weinbauer,  
80 1993; Yamada et al., 2018), can sink faster (Lawrence and Suttle, 2004), and can lead  
81 to preferential grazing by heterotrophic protists (Evans and Wilson, 2008) and/or to  
82 higher growth of grazers (Goode et al., 2019). This process termed “viral shuttle  
83 (Sullivan et al., 2017)” is supported by several field studies that reported association  
84 of viruses with sinking material. Viruses were observed in sinking material in the  
85 North Atlantic Ocean (Proctor and Fuhrman, 1991) and sediment of coastal waters  
86 where algal blooms occur (Lawrence et al., 2002; Tomaru et al., 2007, 2011). In  
87 addition, vertical transport of bacterial viruses between photic and aphotic zones was  
88 observed in the Pacific Ocean (Hurwitz et al., 2015) and in *Tara* Oceans virome data  
89 (Brum et al., 2015). A systematic analysis of large-scale omics data from oligotrophic  
90 oceanic regions revealed a strong positive association between carbon flux and  
91 bacterial dsDNA viruses (*i.e.*, cyanophages), which were previously unrecognized as  
92 possible contributors to the BCP (Guidi et al., 2016). More recently, viral infection of  
93 blooms of the photosynthetic eukaryote *Emiliana huxleyi* in the North Atlantic were

94 found to be accompanied by particle aggregation and greater downward vertical flux  
95 of carbon, with the highest export during the early stage of viral infection (Laber et  
96 al., 2018). These studies raise the question of the overall impact of viruses infecting  
97 eukaryotes on oceanic carbon cycling and export. Given the significant contributions  
98 of eukaryotic plankton to ocean biomass and net production (Hirata et al., 2011; Li,  
99 1995) and their observed predominance over prokaryotes in sinking materials of  
100 Sargasso Sea oligotrophic surface waters (Fawcett et al., 2011; Lomas and Moran,  
101 2011), various lineages of eukaryotic viruses may be responsible for a substantial part  
102 of the variation in carbon export across oceanic regions.

103         If the “viral shunt” and “shuttle” processes function at a global scale and if  
104 these involve specific viruses, we expect to detect a statistical association between  
105 viral community composition and CEE in a large scale omics data. To our knowledge,  
106 such an association has never been investigated. Although this test per se does not  
107 prove that viruses regulate CEE, we consider the association is worth being tested  
108 because such an association is a necessary condition for the global model of viral  
109 shunt and shuttle and, under its absence, we will have to reconsider the model. Deep  
110 sequencing of planktonic community DNA and RNA, as carried out in *Tara Oceans*,  
111 has enabled the identification of marker genes of major viral groups infecting  
112 eukaryotes (Hingamp et al., 2013; Carradec et al., 2018; Culley, 2018). To examine  
113 the association between viral community composition and CEE, we thus used the  
114 comprehensive organismal dataset from the *Tara Oceans* expedition (Carradec et al.,  
115 2018; Sunagawa et al., 2015), as well as related measurements of carbon export  
116 estimated from particle concentrations and size distributions observed *in situ* (Guidi et  
117 al., 2016).

118 In the present study, we identified several hundred marker-gene sequences of  
119 nucleocytoplasmic large DNA viruses (NCLDV) in metagenomes of 0.2–3  $\mu\text{m}$  size  
120 fraction. We also identified RNA and ssDNA viruses in metatranscriptomes of four  
121 eukaryotic size fractions spanning 0.8 to 2,000  $\mu\text{m}$ . The resulting profiles of viral  
122 distributions were compared with an image-based measure of carbon export efficiency  
123 (CEE), which is defined as the ratio of the carbon flux at depth to the carbon flux at  
124 surface.

## 125 **Results and Discussion**

### 126 **Detection of diverse eukaryotic viruses in *Tara Oceans* gene catalogs**

127 We used profile hidden Markov model-based homology searches to identify marker-  
128 gene sequences of eukaryotic viruses in two ocean gene catalogs. These catalogs were  
129 previously constructed from environmental shotgun sequence data of samples  
130 collected during the *Tara Oceans* expedition. The first catalog, the Ocean Microbial  
131 Reference Gene Catalog (OM-RGC), contains 40 million non-redundant genes  
132 predicted from the assemblies of *Tara Oceans* viral and microbial metagenomes  
133 (Sunagawa et al., 2015). We searched this catalog for NCLDV DNA polymerase  
134 family B (PolB) genes, as dsDNA viruses may be present in microbial metagenomes  
135 because large virions ( $> 0.2 \mu\text{m}$ ) have been retained on the filter or because viral  
136 genomes actively replicating or latent within picoeukaryotic cells have been captured.  
137 The second gene catalog, the Marine Atlas of *Tara Oceans* Unigenes (MATOU),  
138 contains 116 million non-redundant genes derived from metatranscriptomes of single-  
139 cell microeukaryotes and small multicellular zooplankton (Carradec et al., 2018). We  
140 searched this catalog for NCLDV PolB genes, RNA-dependent RNA polymerase

141 (RdRP) genes of RNA viruses, and replication-associated protein (Rep) genes of  
142 ssDNA viruses, since transcripts of viruses actively infecting their hosts, as well as  
143 genomes of RNA viruses, have been captured in this catalog.

144 We identified 3,874 NCLDV PolB sequences (3,486 in metagenomes and 388  
145 in metatranscriptomes), 975 RNA virus RdRP sequences, and 299 ssDNA virus Rep  
146 sequences (Table 1). These sequences correspond to operational taxonomic units  
147 (OTUs) at a 95% identity threshold. All except 17 of the NCLDV PolBs from  
148 metagenomes were assigned to the families *Mimiviridae* ( $n = 2,923$ ),  
149 *Phycodnaviridae* ( $n = 348$ ), and *Iridoviridae* ( $n = 198$ ) (Table 1). The larger numbers  
150 of PolB sequences assigned to *Mimiviridae* and *Phycodnaviridae* compared with other  
151 NCLDV families are consistent with a previous observation based on a smaller  
152 dataset (Hingamp et al., 2013). The divergence between these environmental  
153 sequences and reference sequences from known viral genomes was greater in  
154 *Mimiviridae* than in *Phycodnaviridae* (Figure 1a, S1a and S2). Within *Mimiviridae*,  
155 83% of the sequences were most similar to those from algae-infecting *Mimivirus*  
156 relatives. Among the sequences classified in *Phycodnaviridae*, 93% were most similar  
157 to those in *Prasinovirus*, whereas 6% were closest to *Yellowstone lake phycodnavirus*,  
158 which is closely related to *Prasinovirus*. Prasinoviruses are possibly over-represented  
159 in the metagenomes because the 0.2 to 3  $\mu\text{m}$  size fraction selects their picoeukaryotic  
160 hosts. RdRP sequences were assigned mostly to the order *Picornavirales* ( $n = 325$ ),  
161 followed by the families *Partitiviridae* ( $n = 131$ ), *Narnaviridae* ( $n = 95$ ),  
162 *Tombusviridae* ( $n = 45$ ), and *Virgaviridae* ( $n = 33$ ) (Table 1), with most sequences  
163 being distant (30% to 40% amino acid identity) from reference viruses (Figures 1b,  
164 S1b and S3). These results are consistent with previous studies on the diversity of  
165 marine RNA viruses, in which RNA virus sequences were found to correspond to



166 diverse positive-polarity ssRNA and dsRNA viruses distantly related to well-  
167 characterized viruses (Culley, 2018). *Picornavirales* may be over-represented in the  
168 metatranscriptomes because of the polyadenylated RNA selection. The majority ( $n =$   
169 201) of Rep sequences were annotated as *Circoviridae*, known to infect animals,  
170 which is consistent with a previous report (Wang et al., 2018). Only eight were  
171 annotated as plant ssDNA viruses (families *Nanoviridae* and *Gemnaviridae*) (Table  
172 1). Most of these environmental sequences are distant (40% to 50% amino acid  
173 identity) from reference sequences (Figures 1c, S1c and S4). Additional 388 NCLDV  
174 PolBs were detected in the metatranscriptomes. The average cosmopolitanism (number  
175 of samples where an OTU was observed by at least two reads) for PolBs in  
176 metagenomes was 23 samples against 2.9 for metatranscriptome-derived PolB  
177 sequences, 5.5 for Repls, and 5.8 for RdRPs. Within metatranscriptomes, the average  
178 gene-length normalized read counts for PolBs were respectively ten and three times  
179 lower than those of RdRPs and Repls. Therefore, PolBs from metatranscriptomes were  
180 not further used in our study.

## 181 **Composition of eukaryotic viruses can explain the variation of carbon** 182 **export efficiency**

183 Among the PolB, RdRP, and Rep sequences identified in the *Tara* Oceans gene  
184 catalogs, 38%, 18%, and 11% (total = 1,523 sequences), respectively, were present in  
185 at least five samples and had matching carbon export measurement data (Table 1). We  
186 used the relative abundance (defined as the centered log-ratio transformed gene-length  
187 normalized read count) profiles of these 1,523 marker-gene sequences at 59 sampling  
188 sites in the photic zone of 39 *Tara* Oceans stations (Figure 2) to test for association  
189 between their composition and a measure of carbon export efficiency (CEE, see  
190 [Transparent Methods, Figure S5](#)). A partial least squares (PLS) regression model

191 explained 67% (coefficient of determination  $R^2 = 67\%$ ) of the variation in CEE with a  
192 Pearson correlation coefficient of 0.84 between observed and predicted values. This  
193 correlation was confirmed to be statistically significant by permutation test ( $P < 1 \times$   
194  $10^{-4}$ ) (Figure 3a).

195 We also tested for their association with estimates of carbon export flux at 150  
196 meters ( $CE_{150}$ ) and NPP. PLS regressions explained 54% and 64% of the variation in  
197  $CE_{150}$  and NPP with Pearson correlation coefficients between observed and predicted  
198 values of 0.74 (permutation test,  $P < 1 \times 10^{-4}$ ) and 0.80 (permutation test,  $P < 1 \times$   
199  $10^{-4}$ ), respectively (Figure S6). In these three PLS regression models, 83, 86, and 97  
200 viruses were considered to be key predictors (*i.e.*, Variable Importance in the  
201 Projection [VIP] score  $> 2$ ) of CEE,  $CE_{150}$ , and NPP, respectively. PLS models for  
202 NPP and  $CE_{150}$  shared a larger number of predictors (52 viruses) compared to the PLS  
203 models for NPP and CEE (seven viruses) (two proportion Z-test,  $P = 4.14 \times 10^{-12}$ ).  
204 Consistent with this observation,  $CE_{150}$  was correlated with NPP (Pearson's  $r = 0.77$ ;  
205 parametric test,  $P < 1 \times 10^{-12}$ ). This result implies that the magnitude of export in the  
206 analyzed samples was partly constrained by primary productivity. However, CEE was  
207 not correlated with NPP ( $r = 0.16$ ; parametric test,  $P = 0.2$ ) or  $CE_{150}$  ( $r = 0.002$ ;  
208 parametric test,  $P = 0.99$ ). Thus, as expected, primary productivity was not a major  
209 driver for the efficiency of carbon export.

210 The 83 viruses (5% of the viruses included in our analysis) that were  
211 associated with CEE with a VIP score  $> 2$  are considered to be important predictors of  
212 CEE in the PLS regression (Figure 3b, Supplemental Data 1), and these viruses are  
213 hereafter referred to as VIPs (Viruses Important in the Prediction). Fifty-eight VIPs  
214 had positive regression coefficient, and 25 had negative regression coefficient in the  
215 prediction (Figure 3b). Most of the positively associated VIPs showed high relative

216 abundance in the Mediterranean Sea and in the Indian Ocean where CEE tends to be  
217 high compared with other oceanic regions (Figure 4). Among them, 15 (red labels in  
218 Figure 4) also had high relative abundance in samples from other oceanic regions,  
219 showing that these viruses are associated with CEE at a global scale. In contrast,  
220 negatively associated VIPs tend to have higher relative abundance in the Atlantic  
221 Ocean and the Southern Pacific Ocean where CEE is comparatively lower. In the  
222 following sections, we investigate potential hosts of the VIPs in order to interpret the  
223 statistical association between viral community composition and CEE in the light of  
224 previous observations in the literature.

### 225 **Viruses correlated with CEE infect ecologically important hosts**

226 Most of the VIPs (77 of 83) belong to *Mimiviridae* ( $n = 34$  with 25 positive  
227 VIPs and nine negative VIPs), *Phycodnaviridae* ( $n = 24$  with 18 positive VIPs and six  
228 negative VIPs), and ssRNA viruses of the order *Picornavirales* ( $n = 19$  with 13  
229 positive VIPs and six negative VIPs) (Figure 3b, Table S1). All the phycodnavirus  
230 VIPs were most closely related to prasinoviruses infecting Mamiellales, with amino  
231 acid sequence percent identities to reference sequences ranging between 35% and  
232 95%. The six remaining VIPs were two NCLDV of the family *Iridoviridae*  
233 negatively associated with CEE, three RNA viruses (two ssRNA viruses of the family  
234 *Hepeviridae* negatively associated with CEE and one dsRNA virus of the family  
235 *Partitiviridae* positively associated with CEE), and one ssDNA virus of the family  
236 *Circoviridae* positively associated with CEE.

237 Host information may help understand the relationship between these VIPs  
238 and CEE. We performed genomic context analysis for PolB VIPs and phylogeny-  
239 guided network-based host prediction for PolB and RdRP to infer putative virus–host  
240 relationships (see [Transparent Methods](#)).

241 Taxonomic analysis of genes predicted in 10 metagenome-assembled genomes  
242 (MAGs) from the eukaryotic size fractions and 65 genome fragments (contigs)  
243 assembled from the prokaryotic size fraction encoding VIP PolBs further confirmed  
244 their identity as *Mimiviridae* or *Phycodnaviridae* (Figure S7). The size of MAGs  
245 ranged between 30 kbp and 440 kbp with an average of 210 kbp (Table S2). The  
246 presence of genes with high sequence similarities to cellular genes in a viral genome  
247 is suggestive of a virus–host relationship (Monier et al., 2009; Yoshikawa et al.,  
248 2019). Two closely related *Mimiviridae* VIPs, PolB 000079111 (positively associated  
249 with CEE) and PolB 000079078 (negatively associated with CEE), were  
250 phylogenetically close to the pelagophyte virus *Aureococcus anophagefferens virus*  
251 (AaV). One MAG (268 kbp in size) corresponding to PolB 000079111 encoded seven  
252 genes showing high similarities to genes from Pelagophyceae, and another MAG (382  
253 kbp in size), corresponding to PolB 000079078, encoded five genes similar to genes  
254 from Pelagophyceae. All but one of these 12 genes was encoded on a genome  
255 fragment containing genes annotated as viral, including five NCLDV core genes  
256 (Supplemental Data 2), excluding the possibility of contamination in these MAGs.  
257 Two closely related *Phycodnaviridae* VIPs, PolB 001064263 and 010288541, were  
258 positively associated with CEE. Both of these PolBs correspond to a MAG (134 kbp  
259 in size) encoding one gene likely derived from Mamiellales. The genomic fragment  
260 harboring this cellular gene was found to encode 10 genes annotated as viral  
261 (Supplemental Data 2).

262 We conducted a phylogeny-guided, network-based host prediction analysis for  
263 *Mimiviridae*, *Phycodnaviridae*, and *Picornavirales* (Figures S8 and S9). Only a subset  
264 of the VIPs was included in this analysis because we kept the most reliable sequences  
265 to obtain a well-resolved tree topology. Within the *Prasinovirus* clade, which

266 contained thirteen VIPs (nine positive and four negative), seven different eukaryotic  
267 orders were detected as predicted host groups for 10 nodes in the tree. Mamiellales,  
268 the only known host group of prasinoviruses, was detected at eight nodes (five of  
269 them had no parent-to-child relationships), whereas the other six eukaryotic orders  
270 were found at only one node (or two in the case of Eutreptiales) (Figure S8). The  
271 order Mamiellales includes three genera (*Micromonas*, *Ostreococcus*, and  
272 *Bathycoccus*), which are bacterial-sized green microalgae common in coastal and  
273 oceanic environments and are considered to be influential actors in oceanic systems  
274 (Monier et al., 2016). Various prasinoviruses (fourteen with available genome  
275 sequences) have been isolated from the three genera.

276         Within the family *Mimiviridae*, which contains fifteen VIPs (10 positive and  
277 five negative), twelve different orders were predicted as putative host groups (Figure  
278 S8). Collodaria was detected at 15 nodes (two of them had no parent-to-child  
279 relationships), and Prymnesiales at six nodes (three of them had no parent-to-child  
280 relationships), whereas all other orders were present at a maximum of one node each  
281 with no parent-to-child relationships. The nodes enriched for Prymnesiales and  
282 Collodaria fell within a monophyletic clade (marked by a red arrow in Figure S8)  
283 containing four reference haptophyte viruses infecting Prymnesiales and two  
284 reference haptophyte viruses infecting Phaeocystales. Therefore, the environmental  
285 PolB sequences in this *Mimiviridae* clade (including five positive VIPs and one  
286 negative VIP) are predicted to infect Prymnesiales or related haptophytes. The  
287 detection of Collodaria may be the result of indirect associations that reflect a  
288 symbiotic relationship with Prymnesiales, as some acantharians, evolutionarily related  
289 to the Collodaria, are known to host Prymnesiales species (Mars Brisbin et al., 2018).  
290 Known species of Prymnesiales and Phaeocystales have organic scales, except one

291 Prymnesiales species, *Prymnesium neolepis*, which bears siliceous scales (Yoshida et  
292 al., 2006). Some species can form blooms and colonies. Previous studies revealed the  
293 existence of diverse and abundant noncalcifying picohaptophytes in open oceans  
294 (Endo et al., 2018; Liu et al., 2009). Haptophytes as a whole have been estimated to  
295 contribute from 30% to 50% of the total photosynthetic standing stock across the  
296 world ocean (Hirata et al., 2011; Liu et al., 2009). They constitute an important  
297 mixotrophic group in oligotrophic waters (Endo et al., 2018), and mixotrophy is  
298 proposed to increase vertical carbon flux by enabling the uptake of organic forms of  
299 nutrients (Ward and Follows, 2016). Clear host prediction was not made for the other  
300 nine *Mimiviridae* VIPs shown in the phylogenetic tree. Three VIPs (two positive and  
301 one negative) in the tree were relatives of AaV. One negatively associated VIP was a  
302 relative of *Cafeteria roenbergensis virus* infecting a heterotrophic protist. The five  
303 remaining *Mimiviridae* VIPs are very distant from any known *Mimiviridae*.

304         Sixteen *Picornavirales* VIPs (eleven positive and five negative) were included  
305 in the phylogeny-guided, network-based host prediction analysis (Figure 9). Nine  
306 (seven positive and two negative) were grouped within *Dicistroviridae* (known to  
307 infect insects) and may therefore infect marine arthropods such as copepods, the most  
308 ubiquitous and abundant mesozooplankton groups involved in carbon export (Turner,  
309 2015). Three other *Picornavirales* VIPs were placed within a clade containing known  
310 bacillarnaviruses. Two of them (35179764 and 33049404) were positively associated  
311 with CEE and had diatoms of the order Chaetocerotales as a predicted host group. The  
312 third one (107558617) was negatively associated with CEE and distant from other  
313 bacillarnaviruses, and had no host prediction. Diatoms have been globally observed in  
314 the deep sea (Agusti et al., 2015; Leblanc et al., 2018) and identified as important  
315 contributors of the biological carbon pump (Tréguer et al., 2018). One positively

316 associated VIP (32150309) was in a clade containing *Aurantiochytrium single-*  
317 *stranded RNA virus* (AsRNAV), infecting a marine fungoid protist thought to be an  
318 important decomposer (Takao et al., 2005). The last three *Picornavirales* VIPs  
319 (59731273, 49554577, and 36496887) had no predicted host and were too distant  
320 from known *Picornavirales* to speculate about their putative host group.

321 Outside *Picornavirales*, three RNA virus VIPs (two *Hepeviridae*, negatively  
322 associated, and one *Partitiviridae*, positively associated) were identified, for which no  
323 reliable host inferences were made by sequence similarity. Known *Hepeviridae* infect  
324 metazoans, and known *Partitiviridae* infect fungi and plants. The two *Hepeviridae*-  
325 like viruses were most closely related to viruses identified in the transcriptomes of  
326 mollusks (amino acid identities of 48% for 42335229 and 43% for 77677770) (Shi et  
327 al., 2016). The *Partitiviridae*-like VIP (35713768) was most closely related to a  
328 fungal virus, *Penicillium stoloniferum virus S* (49% amino acid identity).

329 One ssDNA virus VIP (38177659) was positively associated with CEE. It was  
330 annotated as a *Circoviridae*, although it groups with other environmental sequences as  
331 an outgroup of known *Circoviridae*. This VIP was connected with copepod, mollusk,  
332 and Collodaria OTUs in the co-occurrence network but no enrichment of predicted  
333 host groups was detected for its clade. *Circoviridae*-like viruses are known to infect  
334 copepods (Dunlap et al., 2013) and have been reported to associate with mollusks  
335 (Dayaram et al., 2015), but none have been reported for Collodaria.

336 Overall, we could infer hosts for 36 VIPs (Tables S3 and S4). Most of the  
337 predicted hosts are known to be ecologically important as primary producers  
338 (Mamiellales, Prymnesiales, Pelagophyceae, and diatoms) or grazers (copepods). Of  
339 these, diatoms and copepods are well known as important contributors to the BCP but  
340 others (*i.e.*, Mamiellales, Prymnesiales, Pelagophyceae) have not been recognized as

341 major contributors to the BCP. Our analysis also revealed that positive and negative  
342 VIPs are not separated in either the viral or host phylogenies.

### 343 **Viruses positively correlated with CEE tend to interact with silicified** 344 **organisms**

345 The phylogeny-guided, network-based host prediction analysis correctly predicted  
346 known virus–host relationships (for viruses infecting Mamiellales, Prymnesiales, and  
347 Chaetocerotales) using our large dataset, despite the reported limitations of these co-  
348 occurrence network-based approaches (Coenen and Weitz, 2018). This result  
349 prompted us to further exploit the species co-occurrence networks (Table S5) to  
350 investigate functional differences between the eukaryotic organisms predicted to  
351 interact with positive VIPs, negative VIPs, and viruses less important for prediction of  
352 CEE (VIP score < 2) (non-VIPs). Positive VIPs had a greater proportion of  
353 connections with silicified eukaryotes ( $Q = 0.001$ ), but not with chloroplast-bearing  
354 eukaryotes ( $Q = 0.16$ ) nor calcifying eukaryotes ( $Q = 1$ ), compared to non-VIPs  
355 (Table S6). No functional differences were observed between negative VIPs and non-  
356 VIPs viruses (Table S6) or positive VIPs (Table S7).

### 357 **Multifarious ways viruses affect the fate of carbon**

358 Our analysis revealed that eukaryotic virus composition was able to predict CEE in  
359 the global sunlit ocean and 83 out of the 1,523 viruses had a high importance in the  
360 predictive model. This association is not a proof that the viruses are the cause of the  
361 variation of CEE. For example, a virus may be found to be associated with CEE if its  
362 host affects CEE regardless of viral infection. This would be the case especially if  
363 latent/persistent viruses are widespread and abundant in phytoplankton (Goic and  
364 Saleh, 2012). Organisms that preferentially grow in marine snow (Bochdansky et al.,



365 2017) may bring associations between viruses infecting those organisms and CEE;  
366 this could be the case for the AsRNAV-related VIP that we identified. Alternatively,  
367 the observed associations between VIPs and CEE may reflect a more direct causal  
368 relationship, which we attempt to explore in light of the large body of literature on the  
369 mechanisms by which viruses impact the fate of carbon in the oceans.

370       Among the 83 VIPs, 58 were positively associated with CEE. Such a positive  
371 association is expected from the “viral shuttle” model, which states that viral activity  
372 could facilitate carbon export to the deep ocean (Fuhrman, 1999; Sullivan et al., 2017;  
373 Weinbauer, 2004) because viral infection can facilitate cell sinking (Lawrence and  
374 Suttle, 2004) and increase the sizes of particles (Peduzzi and Weinbauer, 1993;  
375 Yamada et al., 2018); for instance, a virus may induce secretion of sticky material that  
376 contributes to cell/particle aggregation, such as transparent exopolymeric particles  
377 (TEP) (Nissimov et al., 2018). The data we used to estimate carbon export are based  
378 on the particle size distribution and concentration, and do not convey information  
379 regarding the aggregation status of particles. Therefore, we cannot directly test for a  
380 relation between viruses and aggregation at the sampling sites. Nonetheless, we found  
381 that CEE (*i.e.*,  $CE_{\text{deep}}/CE_{\text{surface}}$ ) increased with the change of particles size from  
382 surface to deep ( $\rho = 0.42$ ,  $P = 8 \times 10^{-9}$ ) (Figure S10). This positive correlation may  
383 reflect an elevated level of aggregation (either enhanced by viral activity or not) in  
384 places where CEE is high, although it could be also due to the presence of large  
385 organisms at depth.

386       Greater aggregate sinking along with higher particulate carbon fluxes was  
387 observed in North Atlantic blooms of *Emiliana huxleyi* that were infected early by  
388 the virus EhV, compared with late-infected blooms (Laber et al., 2018). In the same  
389 bloom, viral infection stage was found to proceed with water column depth (Sheyn et

390 al., 2018). Enhanced TEP production for these same early infected calcifying  
391 populations was observed over a three-day period in deck-board bottle incubations  
392 (Laber et al., 2018). Laboratory observations also exist for enhanced TEP production  
393 and aggregate formation during the early phase of EhV infection of a calcifying *E.*  
394 *huxleyi* strain (Laber et al., 2018; Nissimov et al., 2018). These observations strongly  
395 suggest that infection-induced TEP production in organisms containing dense material  
396 (*e.g.*, calcite scales for *E. huxleyi*) can facilitate carbon export. No EhV-like PolB  
397 sequences were detected in our dataset, which was probably due to sampled areas and  
398 seasons.

399         Laboratory experiments suggest that viruses closely related to positive VIPs,  
400 such as prasinoviruses, have infectious properties that may drive carbon export.  
401 Cultures of *Micromonas pusilla* infected with prasinoviruses showed increased TEP  
402 production compared with non-infected cultures (Lønborg et al., 2013), although it is  
403 not known if this increase leads to aggregation. The hosts of prasinoviruses have been  
404 proposed to contribute to carbon export because they were observed in abyssopelagic  
405 zone at sampling sites dominated by Mamiellales in their surface waters in the  
406 western subtropical North Pacific (Shiozaki et al., 2019). Some prasinoviruses encode  
407 glycosyltransferases (GTs) of the GT2 family. Similar to the a098r gene (GT2) in  
408 *Paramecium bursaria Chlorella virus 1*, the expression of GT2 family members  
409 during infection possibly leads to the production of a dense fibrous hyaluronan  
410 network at the surface of infected cells. Such a network may trigger the aggregation  
411 of host cells, facilitate viral propagation (Van Etten et al., 2017), and increase the cell  
412 wall C:N ratio. We detected one GT2 in a MAG of two *Phycodnaviridae*-like positive  
413 VIPs (000200745 and 002503270) predicted to infect Mamiellales, one in a MAG  
414 corresponding to the putative pelagophyte positive VIP 000079111 related to AaV

415 and six in two MAGs (three each) corresponding to two *Mimiviridae*-like positive  
416 VIPs (000328966 and 001175669). *Phaeocystis globosa virus* (PgV), closely related  
417 to the positive VIP PolB 000912507 (Figure S8), has been linked with increased TEP  
418 production and aggregate formation during the termination of a *Phaeocystis* bloom  
419 (Brussaard et al., 2007). Two closely related bacillarnavirus VIPs were positively  
420 associated with CEE and predicted to infect Chaetocerales. A previous study revealed  
421 an increase in abundance of viruses infecting diatoms of *Chaetoceros* in both the  
422 water columns and the sediments during the bloom of their hosts in a coastal area  
423 (Tomaru et al., 2011), suggesting sinking of cells caused by viruses. Furthermore, the  
424 diatom *Chaetoceros tenuissimus* infected with a DNA virus (CtenDNAV type II) has  
425 been shown to produce higher levels of large-sized particles (50 to 400  $\mu\text{m}$ ) compared  
426 with non-infected cultures (Yamada et al., 2018).

427         The other 25 VIPs were negatively associated with CEE. This association is  
428 compatible with the “viral shunt,” which increases the amount of DOC (Wilhelm and  
429 Suttle, 1999) and reduces the transfer of carbon to higher trophic levels and to the  
430 deep ocean (Fuhrman, 1999; Weitz et al., 2015). Increased DOC has been observed in  
431 culture of Mamiellales lysed by prasinoviruses (Lønborg et al., 2013). Although this  
432 culture-based observation may be difficult to extrapolate to natural conditions, where  
433 the cell concentration and thus the contact rate with viruses are probably lower,  
434 Mamiellales species are known to form blooms during which cell densities may be  
435 comparable with cultures (Zhu et al., 2005). A field study reported that PgV, to which  
436 the negative VIP PolB 000054135 is closely related (Figure S8), can be responsible  
437 for up to 35% of cell lysis per day during bloom of its host (Baudoux et al., 2006),  
438 which is likely accompanied by consequent DOC release. Similarly, the decline of a  
439 bloom of the pelagophyte *Aureococcus anophagefferens* has been associated with

440 active infection by AaV (to which one negative VIP is closely related)  
441 (Moniruzzaman et al., 2017). Among RNA viruses, eight were negative VIPs (six  
442 *Picornavirales* and two *Hepeviridae*). The higher representation of *Picornavirales* in  
443 the viroplankton (Culley, 2018) than within cells (Urayama et al., 2018) suggests that  
444 they are predominantly lytic, although no information exists regarding the effect of  
445 *Picornavirales* on DOC release.

446 It is likely that the “viral shunt” and “viral shuttle” simultaneously affect and  
447 modulate CEE in the global ocean (Zimmerman et al., 2019a). The relative  
448 importance of these two phenomena must fluctuate considerably depending on the  
449 host traits, viral effects on metabolism, and environmental conditions. Reflecting this  
450 complexity, viruses of a same host group could be found to be either positively or  
451 negatively associated with CEE. For example, among prasinoviruses most likely  
452 infecting Mamiellales, 18 were positive VIPs and six were negative VIPs. Two  
453 closely related prasinoviruses (sharing 97.5% genome-wide identity) are known to  
454 exhibit different ecological strategies with notably distinct molecular signatures on the  
455 organic matter released upon infection of the same host (Zimmerman et al., 2019b).  
456 We found that even two very closely related *Mimiviridae* viruses (PolBs 000079111  
457 and 000079078 sharing 94% nucleotide identity over their full gene lengths) most  
458 likely infecting pelagophyte algae were positively and negatively associated with  
459 CEE. Furthermore, it is known that an early-infected *E. huxleyi* system was linked  
460 with both higher aggregation “at surface” and higher remineralization “at deep”  
461 compared to late-infected blooms (Laber et al., 2018). Therefore, the viral effect on  
462 carbon cycle may vary also with depth.

463 Five percent of the tested viruses were associated with CEE in our study.  
464 Similarly, four percent of bacterial virus populations were found to be associated with

465 the magnitude of carbon export at 150 meters (Guidi et al., 2016). These results  
466 suggest that viruses affecting carbon export are rather uncommon. It is plausible that  
467 such viruses affect CEE by infecting organisms that are functionally important  
468 (abundant or keystone species), as we observed in host prediction. The vast majority  
469 (95%) of non-VIPs may not have a significant impact on CEE, because they do not  
470 strongly impact the host population, for instance, by stably coexisting with their hosts.  
471 It is worth noting that experimental studies have reported cultures of algae with  
472 viruses that reach a stable co-existence state after a few generations (Yau et al., 2020).  
473 It is also possible that some of these non-VIPs can impact carbon export but were not  
474 captured in the infection stage affecting the export process. Viruses captured in our  
475 samples can represent active viruses in different infectious stages (early, mid or late)  
476 for metagenomes and metatranscriptomes or at the post-lysis stage for metagenomes.

#### 477 **Potential effects of global climate changes on viral shunt/shuttle and CEE**

478 Increasing evidence suggests that the biological carbon pump is highly dependent on  
479 the planktonic community composition, and as discussed above, viruses represent a  
480 possible key parameter that determines the efficiency of carbon export. In the photic  
481 layer of the oceans, the composition of planktonic communities is strongly affected by  
482 sea surface temperature (Salazar et al., 2019; Sunagawa et al., 2015), and CEE may  
483 therefore be affected by ocean warming. Our result indicated that viruses infecting  
484 small phytoplankton such as Mamiellales and haptophytes are likely associated with  
485 CEE. Interestingly, many studies showed that high temperature and/or CO<sub>2</sub> levels are  
486 associated with an increased contribution of small sized phytoplankton to the total  
487 biomass (Hare et al., 2007; Mousing et al., 2014; Sugie et al., 2020).

488 An increase in CO<sub>2</sub> level in the surface seawater also causes a decrease in pH  
489 (i.e., ocean acidification). Previous studies demonstrated that the decrease in seawater

490 pH negatively affect the growth of calcified and silicified phytoplankton cells (Doney  
491 et al., 2009; Endo et al., 2016; Petrou et al., 2019). The biogenic minerals such as  
492 calcium carbonate and silica act as ballasts in sinking particles (Iversen and Ploug,  
493 2010). Given the statistical association that we detected between the viruses positively  
494 correlated with CEE and the silicified predicted host planktons, the ocean  
495 acidification may decrease the viral shuttle and thus CEE globally in the future.

496         The increased sea surface temperature will decrease the nutrient supply at the  
497 surface of the oligotrophic ocean by preventing the vertical mixing. The decrease in  
498 nutrient availability of surface seawaters possibly diminishes the net primary  
499 production (NPP) and the magnitude of carbon export (corresponded to  $CE_{150}$  in our  
500 study) (Riebesell et al., 2009). Consistently, a downward trend of global  
501 phytoplankton abundance has been observed by satellite (Boyce et al., 2010). In such  
502 a scenario of the global decrease of NPP, the efficiency of export would be an  
503 important factor for a precise estimation of carbon export in the future ocean. In this  
504 regard, the role of marine viruses in the carbon cycle and export should be further  
505 investigated and eventually be integrated into prospective models for the climate  
506 change.

## 507 **Conclusions**

508 Eukaryotic virus community composition was able to predict CEE at 59 sampling  
509 sites in the photic zone of the world ocean. This statistical association was detected  
510 based on a large omics dataset collected throughout the oceans and processed with  
511 standardized protocols. The predictability of CEE by viral composition is consistent  
512 with the hypothesis that “viral shuttle” and “shunt” are functioning at a global scale.  
513 Among 83 viruses with a high importance in the prediction of CEE, 58 viruses were  
514 positively and 25 negatively correlated with carbon export efficiency. Most of these

515 viruses belong to *Prasinovirus*, *Mimiviridae*, and *Picornavirales* and are either new to  
516 science or with no known roles in carbon export efficiency. Thirty-six of these  
517 “select” viruses were predicted to infect ecologically important hosts such as green  
518 algae of the order Mamiellales, haptophytes, diatoms, and copepods. Positively  
519 associated viruses had more predicted interactions with silicified eukaryotes than non-  
520 associated viruses did. Overall, these results imply that the effect of viruses on the  
521 “shuttle” and “shunt” processes could be dependent on viral hosts and ecosystem  
522 dynamics.

## 523 **Limitations of the study**

524 The observed statistical associations between viral compositions and examined  
525 parameters (*i.e.*, CEE, CE and NPP) do not convey the information about the direction  
526 of their potential causality relationships, and they could even result from indirect  
527 relationships as discussed above. Certain groups of viruses detected in samples may  
528 be over- or under-represented because of the technical limitations in size  
529 fractionation, DNA/RNA extraction and sequencing.

## 530 **Resource Availability**

### 531 **Lead Contact**

532 Further information and requests for resources and reagents should be directed to and  
533 will be fulfilled by Lead Contact, Hiroyuki Ogata ([ogata@kuicr.kyoto-u.ac.jp](mailto:ogata@kuicr.kyoto-u.ac.jp)).

### 534 **Materials Availability**

535 This study did not generate unique reagent.

## 536 **Data and Code Availability**

537 The authors declare that the data supporting the findings of this study are available

538 within the paper and its supplemental files, as well as at the GenomeNet FTP:

539 [ftp://ftp.genome.jp/pub/db/community/tara/Cpump/Supplementary\\_material/](ftp://ftp.genome.jp/pub/db/community/tara/Cpump/Supplementary_material/).

540 Our custom R script used to test for association between viruses and environmental

541 variables (CEE, CE<sub>150</sub> and NPP) is available along with input data at the GenomeNet

542 FTP:

543 [ftp://ftp.genome.jp/pub/db/community/tara/Cpump/Supplementary\\_material/PLSreg/](ftp://ftp.genome.jp/pub/db/community/tara/Cpump/Supplementary_material/PLSreg/).

544 The Taxon Interaction Mapper (TIM) tool developed for this study and used for virus

545 host prediction is available at <https://github.com/RomainBlancMathieu/TIM>.

## 546 **Supplemental Files**

547 • Supplemental\_Information.pdf: supplemental figures and tables, and transparent

548 methods

549 • Supplemental\_Data\_1\_2.xlsx

## 550 **Acknowledgements**

551 We thank the *Tara* Oceans consortium, the projects Oceanomics and France

552 Genomique (grants ANR-11-BTBR-0008 and ANR-10-INBS-09), and the people and

553 sponsors who supported the *Tara* Oceans Expedition (<http://www.embl.de/tara->

554 oceans/) for making the data accessible. This is contribution number XXX of the *Tara*

555 Oceans Expedition 2009–2013. Computational time was provided by the

556 SuperComputer System, Institute for Chemical Research, Kyoto University. We thank

557 Barbara Goodson, Ph.D., and Sara J. Mason, M.Sc., from Edanz Group (<https://en->

558 author-services.edanzgroup.com/) for editing a draft of this manuscript. This work



559 was supported by JSPS/KAKENHI (Nos. 26430184, 18H02279, and 19H05667 to  
560 H.O. and Nos. 19K15895 and 19H04263 to H.E.), Scientific Research on Innovative  
561 Areas from the Ministry of Education, Culture, Science, Sports and Technology  
562 (MEXT) of Japan (Nos. 16H06429, 16K21723, and 16H06437 to H.O.), the  
563 Collaborative Research Program of the Institute for Chemical Research, Kyoto  
564 University (2019-29 to S.C.), the Future Development Funding Program of the Kyoto  
565 University Research Coordination Alliance (to R.B.M.), the ICR-KU International  
566 Short-term Exchange Program for Young Researchers (to S.C.), and the Research  
567 Unit for Development of Global Sustainability (to H.O. and T.O.D.).

## 568 **Author contributions**

569 H.O. and R.B.M. conceived the study. R.B.M and H.K. performed most of the  
570 analyses. H.E. and L.G. designed carbon export analysis. R.H.V and S.C. performed  
571 network analysis. N.H. and C.d.V. analyzed eukaryotic sequences. T.O.D., M.G., P.F.  
572 and O.J. analyzed viral MAGs. C.H.N. and H.M. contributed to statistical analysis.  
573 M.B.S. and C.A.S. contributed to interpretations. All authors edited and approved the  
574 final version of the manuscript.

## 575 **Declaration of Interests**

576 The authors declare no competing interests.

## 577 **References**

578 Agusti, S., González-Gordillo, J.I., Vaqué, D., Estrada, M., Cerezo, M.I., Salazar, G.,  
579 Gasol, J.M., and Duarte, C.M. (2015). Ubiquitous healthy diatoms in the deep sea  
580 confirm deep carbon injection by the biological pump. *Nat. Commun.* 6, 7608.

- 581 Baudoux, A., Noordeloos, A., Veldhuis, M., and Brussaard, C. (2006). Virally  
582 induced mortality of *Phaeocystis globosa* during two spring blooms in temperate  
583 coastal waters. *Aquat. Microb. Ecol.* *44*, 207–217.
- 584 Bochdansky, A.B., Clouse, M.A., and Herndl, G.J. (2017). Eukaryotic microbes,  
585 principally fungi and labyrinthulomycetes, dominate biomass on bathypelagic marine  
586 snow. *ISME J.* *11*, 362–373.
- 587 Boyce, D.G., Lewis, M.R., and Worm, B. (2010). Global phytoplankton decline over  
588 the past century. *Nature* *466*, 591–596.
- 589 Brum, J.R., Ignacio-Espinoza, J.C., Roux, S., Doucier, G., Acinas, S.G., Alberti, A.,  
590 Chaffron, S., Cruaud, C., Vargas, C. de, Gasol, J.M., et al. (2015). Patterns and  
591 ecological drivers of ocean viral communities. *Science* *348*, 1261498.
- 592 Brussaard, C.P.D., Bratbak, G., Baudoux, A.-C., and Ruardij, P. (2007). *Phaeocystis*  
593 and its interaction with viruses. *Biogeochemistry* *83*, 201–215.
- 594 Buesseler, K.O., and Boyd, P.W. (2009). Shedding light on processes that control  
595 particle export and flux attenuation in the twilight zone of the open ocean. *Limnol.*  
596 *Oceanogr.* *54*, 1210–1232.
- 597 Carradec, Q., Pelletier, E., Silva, C.D., Alberti, A., Seeleuthner, Y., Blanc-Mathieu,  
598 R., Lima-Mendez, G., Rocha, F., Tirichine, L., Labadie, K., et al. (2018). A global  
599 ocean atlas of eukaryotic genes. *Nat. Commun.* *9*, 373.
- 600 Coenen, A.R., and Weitz, J.S. (2018). Limitations of Correlation-Based Inference in  
601 Complex Virus-Microbe Communities. *MSystems* *3*, e00084-18.
- 602 Culley, A. (2018). New insight into the RNA aquatic virosphere via viromics. *Virus*  
603 *Res.* *244*, 84–89.
- 604 Dayaram, A., Goldstien, S., Argüello-Astorga, G.R., Zawar-Reza, P., Gomez, C.,  
605 Harding, J.S., and Varsani, A. (2015). Diverse small circular DNA viruses circulating  
606 amongst estuarine molluscs. *Infect. Genet. Evol. J. Mol. Epidemiol. Evol. Genet.*  
607 *Infect. Dis.* *31*, 284–295.
- 608 Doney, S.C., Fabry, V.J., Feely, R.A., and Kleypas, J.A. (2009). Ocean Acidification:  
609 The Other CO<sub>2</sub> Problem. *Annu. Rev. Mar. Sci.* *1*, 169–192.
- 610 Dunlap, D.S., Ng, T.F.F., Rosario, K., Barbosa, J.G., Greco, A.M., Breitbart, M., and  
611 Hewson, I. (2013). Molecular and microscopic evidence of viruses in marine  
612 copepods. *Proc. Natl. Acad. Sci.* *110*, 1375–1380.
- 613 Endo, H., Sugie, K., Yoshimura, T., and Suzuki, K. (2016). Response of Spring  
614 Diatoms to CO<sub>2</sub> Availability in the Western North Pacific as Determined by Next-  
615 Generation Sequencing. *PLOS ONE* *11*, e0154291.
- 616 Endo, H., Ogata, H., and Suzuki, K. (2018). Contrasting biogeography and diversity  
617 patterns between diatoms and haptophytes in the central Pacific Ocean. *Sci. Rep.* *8*,  
618 10916.

- 619 Evans, C., and Wilson, W.H. (2008). Preferential grazing of *Oxyrrhis marina* on virus  
620 infected *Emiliana huxleyi*. *Limnol. Oceanogr.* 53, 2035–2040.
- 621 Fawcett, S.E., Lomas, M.W., Casey, J.R., Ward, B.B., and Sigman, D.M. (2011).  
622 Assimilation of upwelled nitrate by small eukaryotes in the Sargasso Sea. *Nat.*  
623 *Geosci.* 4, 717–722.
- 624 Fuhrman, J.A. (1999). Marine viruses and their biogeochemical and ecological  
625 effects. *Nature* 399, 541–548.
- 626 Gobler, C.J., Hutchins, D.A., Fisher, N.S., Coper, E.M., and Sañudo- Wilhelmy,  
627 S.A. (1997). Release and bioavailability of C, N, P Se, and Fe following viral lysis of  
628 a marine chrysophyte. *Limnol. Oceanogr.* 42, 1492–1504.
- 629 Goic, B., and Saleh, M.-C. (2012). Living with the enemy: viral persistent infections  
630 from a friendly viewpoint. *Curr. Opin. Microbiol.* 15, 531–537.
- 631 Goode, A.G., Fields, D.M., Archer, S.D., and Martínez, J.M. (2019). Physiological  
632 responses of *Oxyrrhis marina* to a diet of virally infected *Emiliana huxleyi*. *PeerJ* 7,  
633 e6722.
- 634 Guidi, L., Chaffron, S., Bittner, L., Eveillard, D., Larhlimi, A., Roux, S., Darzi, Y.,  
635 Audic, S., Berline, L., Brum, J.R., et al. (2016). Plankton networks driving carbon  
636 export in the oligotrophic ocean. *Nature* 532, 465.
- 637 Hare, C., Leblanc, K., DiTullio, G., Kudela, R., Zhang, Y., Lee, P., Riseman, S., and  
638 Hutchins, D. (2007). Consequences of increased temperature and CO<sub>2</sub> for  
639 phytoplankton community structure in the Bering Sea. *Mar. Ecol. Prog. Ser.* 352, 9–  
640 16.
- 641 Hingamp, P., Grimsley, N., Acinas, S.G., Clerissi, C., Subirana, L., Poulain, J.,  
642 Ferrera, I., Sarmiento, H., Villar, E., Lima-Mendez, G., et al. (2013). Exploring  
643 nucleo-cytoplasmic large DNA viruses in Tara Oceans microbial metagenomes. *ISME*  
644 *J.* 7, 1678–1695.
- 645 Hirata, T., Hardman-Mountford, N.J., Brewin, R.J.W., Aiken, J., Barlow, R., Suzuki,  
646 K., Isada, T., Howell, E., Hashioka, T., Noguchi-Aita, M., et al. (2011). Synoptic  
647 relationships between surface Chlorophyll-*a* and diagnostic pigments specific to  
648 phytoplankton functional types. *Biogeosciences* 8, 311–327.
- 649 Hurwitz, B.L., Brum, J.R., and Sullivan, M.B. (2015). Depth-stratified functional and  
650 taxonomic niche specialization in the “core” and “flexible” Pacific Ocean Virome.  
651 *ISME J.* 9, 472–484.
- 652 Iversen, M.H., and Ploug, H. (2010). Ballast minerals and the sinking carbon flux in  
653 the ocean: carbon-specific respiration rates and sinking velocity of marine snow  
654 aggregates. *Biogeosciences* 7, 2613–2624.
- 655 Karl, D.M., Church, M.J., Dore, J.E., Letelier, R.M., and Mahaffey, C. (2012).  
656 Predictable and efficient carbon sequestration in the North Pacific Ocean supported  
657 by symbiotic nitrogen fixation. *Proc. Natl. Acad. Sci.* 109, 1842–1849.

- 658 Klaas, C., and Archer, D.E. (2002). Association of sinking organic matter with  
659 various types of mineral ballast in the deep sea: Implications for the rain ratio. *Glob.*  
660 *Biogeochem. Cycles* *16*, 63-1-63-14.
- 661 Laber, C.P., Hunter, J.E., Carvalho, F., Collins, J.R., Hunter, E.J., Schieler, B.M.,  
662 Boss, E., More, K., Frada, M., Thamatrakoln, K., et al. (2018). Coccolithovirus  
663 facilitation of carbon export in the North Atlantic. *Nat. Microbiol.* *3*, 537–547.
- 664 Lawrence, J.E., and Suttle, C.A. (2004). Effect of viral infection on sinking rates of  
665 *Heterosigma akashiwo* and its implications for bloom termination. *Aquat. Microb.*  
666 *Ecol.* *37*, 1–7.
- 667 Lawrence, J.E., Chan, A.M., and Suttle, C.A. (2002). Viruses causing lysis of the  
668 toxic bloom-forming alga *Heterosigma akashiwo* (Raphidophyceae) are widespread in  
669 coastal sediments of British Columbia, Canada. *Limnol. Oceanogr.* *47*, 545–550.
- 670 Leblanc, K., Quéguiner, B., Diaz, F., Cornet, V., Michel-Rodriguez, M., Durrieu de  
671 Madron, X., Bowler, C., Malviya, S., Thyssen, M., Grégori, G., et al. (2018).  
672 Nanoplanktonic diatoms are globally overlooked but play a role in spring blooms and  
673 carbon export. *Nat. Commun.* *9*, 953.
- 674 Li, W. (1995). Composition of Ultraphytoplankton in the Central North-Atlantic. *Mar.*  
675 *Ecol. Prog. Ser.* *122*, 1–8.
- 676 Liu, H., Probert, I., Uitz, J., Claustre, H., Aris-Brosou, S., Frada, M., Not, F., and de  
677 Vargas, C. (2009). Extreme diversity in noncalcifying haptophytes explains a major  
678 pigment paradox in open oceans. *Proc. Natl. Acad. Sci. U. S. A.* *106*, 12803–12808.
- 679 Lomas, M.W., and Moran, S.B. (2011). Evidence for aggregation and export of  
680 cyanobacteria and nano-eukaryotes from the Sargasso Sea euphotic zone.  
681 *Biogeosciences* *8*, 203–216.
- 682 Lønborg, C., Middelboe, M., and Brussaard, C.P.D. (2013). Viral lysis of  
683 *Micromonas pusilla*: impacts on dissolved organic matter production and  
684 composition. *Biogeochemistry* *116*, 231–240.
- 685 Mars Brisbin, M., Mesrop, L.Y., Grossmann, M.M., and Mitarai, S. (2018). Intra-host  
686 Symbiont Diversity and Extended Symbiont Maintenance in Photosymbiotic  
687 *Acantharea* (Clade F). *Front. Microbiol.* *9*.
- 688 Monier, A., Pagarete, A., de Vargas, C., Allen, M.J., Read, B., Claverie, J.-M., and  
689 Ogata, H. (2009). Horizontal gene transfer of an entire metabolic pathway between a  
690 eukaryotic alga and its DNA virus. *Genome Res.* *19*, 1441–1449.
- 691 Monier, A., Worden, A.Z., and Richards, T.A. (2016). Phylogenetic diversity and  
692 biogeography of the Mamiellophyceae lineage of eukaryotic phytoplankton across the  
693 oceans. *Environ. Microbiol. Rep.* *8*, 461–469.
- 694 Moniruzzaman, M., Wurch, L.L., Alexander, H., Dyhrman, S.T., Gobler, C.J., and  
695 Wilhelm, S.W. (2017). Virus-host relationships of marine single-celled eukaryotes  
696 resolved from metatranscriptomics. *Nat. Commun.* *8*, 16054.

- 697 Mousing, E., Ellegaard, M., and Richardson, K. (2014). Global patterns in  
698 phytoplankton community size structure—evidence for a direct temperature effect.  
699 *Mar. Ecol. Prog. Ser.* *497*, 25–38.
- 700 Nissimov, J.I., Vandzura, R., Johns, C.T., Natale, F., Haramaty, L., and Bidle, K.D.  
701 (2018). Dynamics of transparent exopolymer particle production and aggregation  
702 during viral infection of the coccolithophore, *Emiliana huxleyi*. *Environ. Microbiol.*  
703 *20*, 2880–2897.
- 704 Peduzzi, P., and Weinbauer, M.G. (1993). Effect of concentrating the virus-rich 2-  
705 2nm size fraction of seawater on the formation of algal flocs (marine snow). *Limnol.*  
706 *Oceanogr.* *38*, 1562–1565.
- 707 Petrou, K., Baker, K.G., Nielsen, D.A., Hancock, A.M., Schulz, K.G., and Davidson,  
708 A.T. (2019). Acidification diminishes diatom silica production in the Southern Ocean.  
709 *Nat. Clim. Change* *9*, 781–786.
- 710 Proctor, L.M., and Fuhrman, J.A. (1991). Roles of viral infection in organic particle  
711 flux. *Mar. Ecol. Prog. Ser.* *69*, 133–142.
- 712 Riebesell, U., Kortzinger, A., and Oschlies, A. (2009). Sensitivities of marine carbon  
713 fluxes to ocean change. *Proc. Natl. Acad. Sci.* *106*, 20602–20609.
- 714 Salazar, G., Paoli, L., Alberti, A., Huerta-Cepas, J., Ruscheweyh, H.-J., Cuenca, M.,  
715 Field, C.M., Coelho, L.P., Cruaud, C., Engelen, S., et al. (2019). Gene Expression  
716 Changes and Community Turnover Differentially Shape the Global Ocean  
717 Metatranscriptome. *Cell* *179*, 1068-1083.e21.
- 718 Sheyn, U., Rosenwasser, S., Lehahn, Y., Barak-Gavish, N., Rotkopf, R., Bidle, K.D.,  
719 Koren, I., Schatz, D., and Vardi, A. (2018). Expression profiling of host and virus  
720 during a coccolithophore bloom provides insights into the role of viral infection in  
721 promoting carbon export. *ISME J.* *1*.
- 722 Shi, M., Lin, X.-D., Tian, J.-H., Chen, L.-J., Chen, X., Li, C.-X., Qin, X.-C., Li, J.,  
723 Cao, J.-P., Eden, J.-S., et al. (2016). Redefining the invertebrate RNA virosphere.  
724 *Nature* *540*, 539–543.
- 725 Shiozaki, T., Hirose, Y., Hamasaki, K., Kaneko, R., Ishikawa, K., and Harada, N.  
726 (2019). Eukaryotic Phytoplankton Contributing to a Seasonal Bloom and Carbon  
727 Export Revealed by Tracking Sequence Variants in the Western North Pacific. *Front.*  
728 *Microbiol.* *10*.
- 729 Sugie, K., Fujiwara, A., Nishino, S., Kameyama, S., and Harada, N. (2020). Impacts  
730 of Temperature, CO<sub>2</sub>, and Salinity on Phytoplankton Community Composition in the  
731 Western Arctic Ocean. *Front. Mar. Sci.* *6*, 821.
- 732 Sullivan, M.B., Weitz, J.S., and Wilhelm, S. (2017). Viral ecology comes of age.  
733 *Environ. Microbiol. Rep.* *9*, 33–35.
- 734 Sunagawa, S., Coelho, L.P., Chaffron, S., Kultima, J.R., Labadie, K., Salazar, G.,  
735 Djahanschiri, B., Zeller, G., Mende, D.R., Alberti, A., et al. (2015). Ocean plankton.  
736 Structure and function of the global ocean microbiome. *Science* *348*, 1261359.

- 737 Suttle, C.A. (2007). Marine viruses--major players in the global ecosystem. *Nat. Rev.*  
738 *Microbiol.* *5*, 801–812.
- 739 Takao, Y., Nagasaki, K., Mise, K., Okuno, T., and Honda, D. (2005). Isolation and  
740 characterization of a novel single-stranded RNA Virus infectious to a marine fungoid  
741 protist, *Schizochytrium* sp. (Thraustochytriaceae, Labyrinthulea). *Appl. Environ.*  
742 *Microbiol.* *71*, 4516–4522.
- 743 Tomaru, Y., Hata, N., Masuda, T., Tsuji, M., Igata, K., Masuda, Y., Yamatogi, T.,  
744 Sakaguchi, M., and Nagasaki, K. (2007). Ecological dynamics of the bivalve-killing  
745 dinoflagellate *Heterocapsa circularisquama* and its infectious viruses in different  
746 locations of western Japan. *Environ. Microbiol.* *9*, 1376–1383.
- 747 Tomaru, Y., Fujii, N., Oda, S., Toyoda, K., and Nagasaki, K. (2011). Dynamics of  
748 diatom viruses on the western coast of Japan. *Aquat. Microb. Ecol.* *63*, 223–230.
- 749 Tréguer, P., Bowler, C., Moriceau, B., Dutkiewicz, S., Gehlen, M., Aumont, O.,  
750 Bittner, L., Dugdale, R., Finkel, Z., Iudicone, D., et al. (2018). Influence of diatom  
751 diversity on the ocean biological carbon pump. *Nat. Geosci.* *11*, 27–37.
- 752 Turner, J.T. (2015). Zooplankton fecal pellets, marine snow, phytodetritus and the  
753 ocean's biological pump. *Prog. Oceanogr.* *130*, 205–248.
- 754 Urayama, S., Takaki, Y., Nishi, S., Yoshida- Takashima, Y., Deguchi, S., Takai, K.,  
755 and Nunoura, T. (2018). Unveiling the RNA virosphere associated with marine  
756 microorganisms. *Mol. Ecol. Resour.* *18*, 1444–1455.
- 757 Van Etten, J.L., Agarkova, I., Dunigan, D.D., Tonetti, M., De Castro, C., and Duncan,  
758 G.A. (2017). Chloroviruses Have a Sweet Tooth. *Viruses* *9*.
- 759 Wang, H., Wu, S., Li, K., Pan, Y., Yan, S., and Wang, Y. (2018). Metagenomic  
760 analysis of ssDNA viruses in surface seawater of Yangshan Deep-Water Harbor,  
761 Shanghai, China. *Mar. Genomics* *41*, 50–53.
- 762 Ward, B.A., and Follows, M.J. (2016). Marine mixotrophy increases trophic transfer  
763 efficiency, mean organism size, and vertical carbon flux. *Proc. Natl. Acad. Sci.* *113*,  
764 2958–2963.
- 765 Weinbauer, M.G. (2004). Ecology of prokaryotic viruses. *FEMS Microbiol. Rev.* *28*,  
766 127–181.
- 767 Weitz, J.S., Stock, C.A., Wilhelm, S.W., Bourouiba, L., Coleman, M.L., Buchan, A.,  
768 Follows, M.J., Fuhrman, J.A., Jover, L.F., Lennon, J.T., et al. (2015). A multitrophic  
769 model to quantify the effects of marine viruses on microbial food webs and ecosystem  
770 processes. *ISME J.* *9*, 1352–1364.
- 771 Wilhelm, S.W., and Suttle, C.A. (1999). Viruses and Nutrient Cycles in the  
772 SeaViruses play critical roles in the structure and function of aquatic food webs.  
773 *BioScience* *49*, 781–788.
- 774 Yamada, Y., Tomaru, Y., Fukuda, H., and Nagata, T. (2018). Aggregate Formation  
775 During the Viral Lysis of a Marine Diatom. *Front. Mar. Sci.* *5*.

- 776 Yau, S., Krasovec, M., Benites, L.F., Rombauts, S., Groussin, M., Vancaester, E.,  
777 Aury, J.-M., Derelle, E., Desdevises, Y., Escande, M.-L., et al. (2020). Virus-host  
778 coexistence in phytoplankton through the genomic lens. *Sci. Adv.* 6, eaay2587.
- 779 Yoshida, M., Noël, M.-H., Nakayama, T., Naganuma, T., and Inouye, I. (2006). A  
780 haptophyte bearing siliceous scales: ultrastructure and phylogenetic position of  
781 *Hyalolithus neolepis* gen. et sp. nov. (Prymnesiophyceae, Haptophyta). *Protist* 157,  
782 213–234.
- 783 Yoshikawa, G., Blanc-Mathieu, R., Song, C., Kayama, Y., Mochizuki, T., Murata, K.,  
784 Ogata, H., and Takemura, M. (2019). Medusavirus, a novel large DNA virus  
785 discovered from hot spring water. *J. Virol.* JVI.02130-18.
- 786 Zhang, C., Dang, H., Azam, F., Benner, R., Legendre, L., Passow, U., Polimene, L.,  
787 Robinson, C., Suttle, C.A., and Jiao, N. (2018). Evolving paradigms in biological  
788 carbon cycling in the ocean. *Natl. Sci. Rev.* 5, 481–499.
- 789 Zhu, F., Massana, R., Not, F., Marie, D., and Vaulot, D. (2005). Mapping of  
790 picoeucaryotes in marine ecosystems with quantitative PCR of the 18S rRNA gene.  
791 *FEMS Microbiol. Ecol.* 52, 79–92.
- 792 Zimmerman, A.E., Howard-Varona, C., Needham, D.M., John, S.G., Worden, A.Z.,  
793 Sullivan, M.B., Waldbauer, J.R., and Coleman, M.L. (2019a). Metabolic and  
794 biogeochemical consequences of viral infection in aquatic ecosystems. *Nat. Rev.*  
795 *Microbiol.* 1–14.
- 796 Zimmerman, A.E., Bachy, C., Ma, X., Roux, S., Jang, H.B., Sullivan, M.B.,  
797 Waldbauer, J.R., and Worden, A.Z. (2019b). Closely related viruses of the marine  
798 picoeukaryotic alga *Ostreococcus lucimarinus* exhibit different ecological strategies.  
799 *Environ. Microbiol.* 21, 2148–2170.
- 800
- 801

## 802 **Figure legends**

803 **Figure 1: Viruses of eukaryotic plankton identified in *Tara* Oceans samples are**  
804 **distantly related to characterized viruses.** Unrooted maximum likelihood  
805 phylogenetic trees containing environmental (black) and reference (red) viral  
806 sequences for NCLDV DNA polymerase family B (**a**), RNA virus RNA-dependent  
807 RNA polymerase (**b**), and ssDNA virus replication-associated protein (**c**). A  
808 rectangular representation of these trees with branch support values is provided in  
809 Figure S2–S4.

810

811 **Figure 2: Carbon export efficiency and relative marker-gene occurrence of**  
812 **eukaryotic plankton viruses along the sampling route. a** Carbon export efficiency  
813 estimated at 39 *Tara* Oceans stations where surface and DCM layers were sampled  
814 for prokaryote-enriched metagenomes and eukaryotic metatranscriptomes. **b and c**  
815 Relative marker-gene occurrence of major groups of viruses of eukaryotic plankton  
816 for NCLDVs in metagenomes (**b**) and for RNA and ssDNA viruses in  
817 metatranscriptomes (**c**) at 59 sampling sites.

818

819 **Figure 3: Relative abundance of eukaryotic plankton viruses associated with**  
820 **carbon export efficiency in the global ocean. a** Bivariate plot between predicted and  
821 observed values in a leave-one-out cross-validation test for carbon export efficiency.  
822 The PLS regression model was constructed using occurrence profiles of 1,523  
823 marker-gene sequences (1,309 PolBs, 180 RdRPs and 34 Repls) derived from  
824 environmental samples.  $r$ , Pearson correlation coefficient;  $R^2$ , the coefficient of  
825 determination between measured response values and predicted response values.  $R^2$ ,  
826 which was calculated as  $1 - \text{SSE}/\text{SST}$  (sum of squares due to error and total)



827 measures how successful the fit is in explaining the variance of the response values.  
828 The significance of the association was assessed using a permutation test ( $n = 10,000$ )  
829 (grey histogram in **a**). The red diagonal line shows the theoretical curve for perfect  
830 prediction. **b** Pearson correlation coefficients between CEE and occurrence profiles of  
831 83 viruses that have VIP scores  $> 2$  (VIPs) with the first two components in the PLS  
832 regression model using all samples. PLS components 1 and 2 explained 83% and 11%  
833 of the variance of CEE, respectively. Fifty-eight VIPs had positive regression  
834 coefficients in the model (shown with circles), and 25 had negative regression  
835 coefficients (shown with triangles).

836

837 **Figure 4: Biogeography of viruses associated with carbon export efficiency.** The  
838 upper panel shows carbon export efficiency ( $CEE = CE_{\text{deep}}/CE_{\text{surface}}$ ) for 59 sampling  
839 sites. The bottom panel is a map reflecting relative abundances, expressed as centered  
840 log-ratio transformed, gene-length normalized read counts of viruses positively and  
841 negatively associated with CEE that have VIP scores  $> 2$  (VIPs). MS, Mediterranean  
842 Sea; IO, Indian Ocean; SAO, South Atlantic Ocean; SPO, South Pacific Ocean; NPO,  
843 North Pacific Ocean; NAO, North Atlantic Ocean. The bottom horizontal axis is  
844 labeled with *Tara* Oceans station numbers, sampling depth (SRF, surface; DCM, deep  
845 chlorophyll maximum), and abbreviations of biogeographic provinces. Viruses  
846 labeled in red correspond to positive VIPs that tend more represented in more than  
847 one biogeographic province.

848

849 **Tables**

850 **Table 1: Taxonomic breakdown of viral marker genes**

Viruses		Identified	Used in PLS regression*
NCLDVs	Mimiviridae	2,923	1,148
	Phycodnaviridae	348	99
	Iridoviridae	198	59
	Other NCLDVs **	17	3
	Total	3,486	1,309
RNA viruses	Picornavirales (ssRNA+)	325	80
	Partitiviridae (dsRNA)	131	22
	Narnaviridae (ssRNA+)	95	6
	Other families	289	53
	Unclassified	78	9
	RNA viruses	57	10
Total	975	180	
ssDNA viruses	Circoviridae	201	22
	Geminiviridae	4	0
	Nanoviridae	4	0
	Unclassified	39	2
	ssDNA viruses	51	10
Total	299	34	
All		4,760	1,523

851 \* The marker genes had to occurred in at least five samples and harbor a Spearman correlation  
 852 coefficient > |0.2| with carbon export efficiency.

853 \*\* There was no unclassified NCLDV.

854

Figure 1-4  
Figure 1

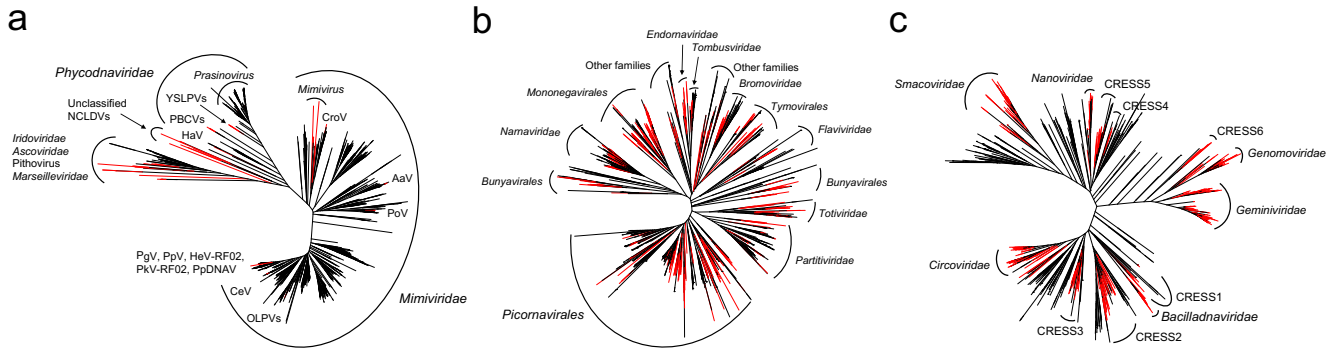


Figure 2

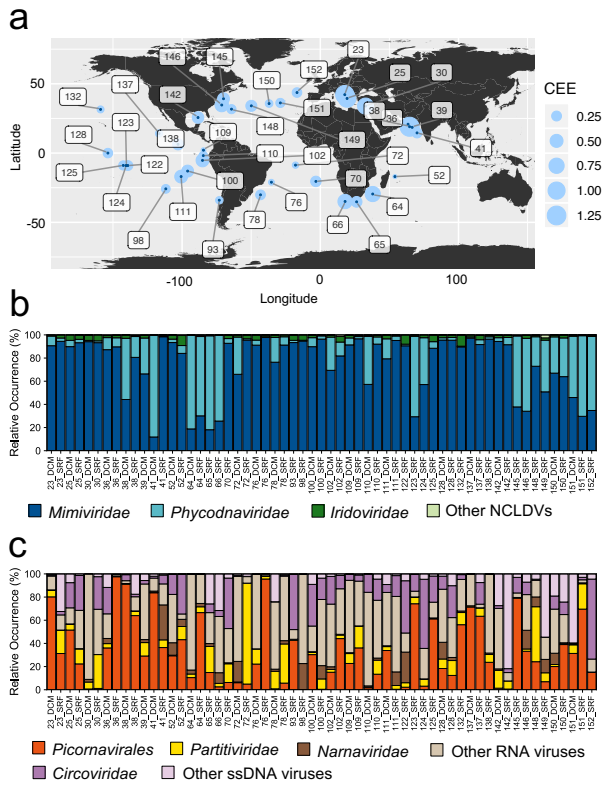
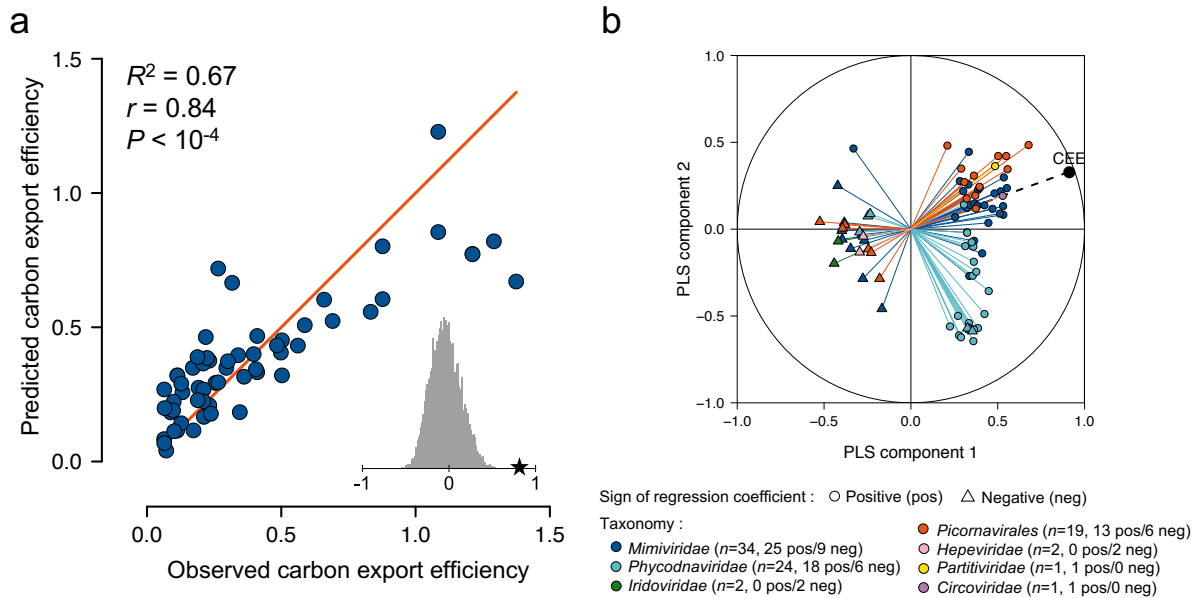
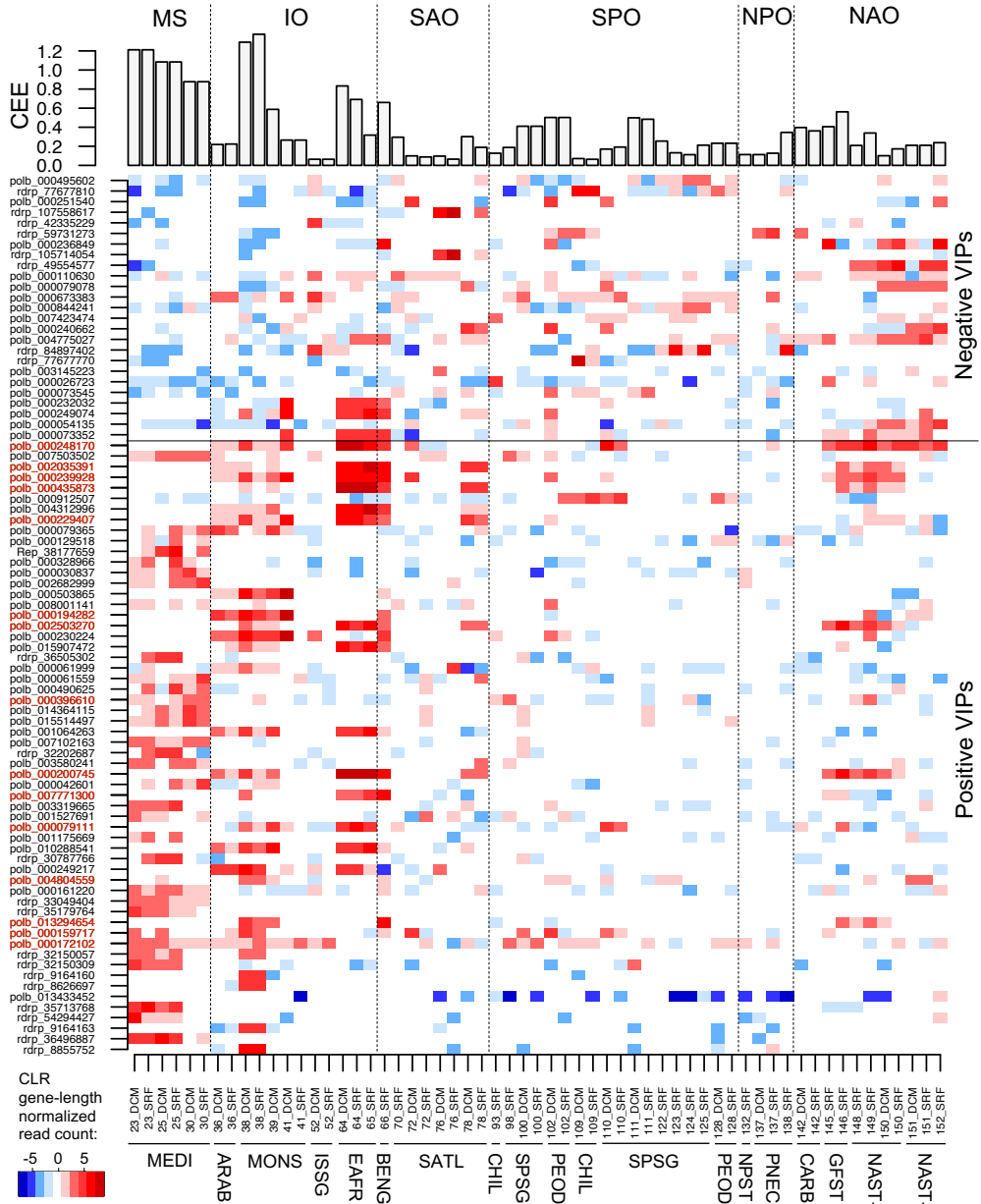


Figure 3



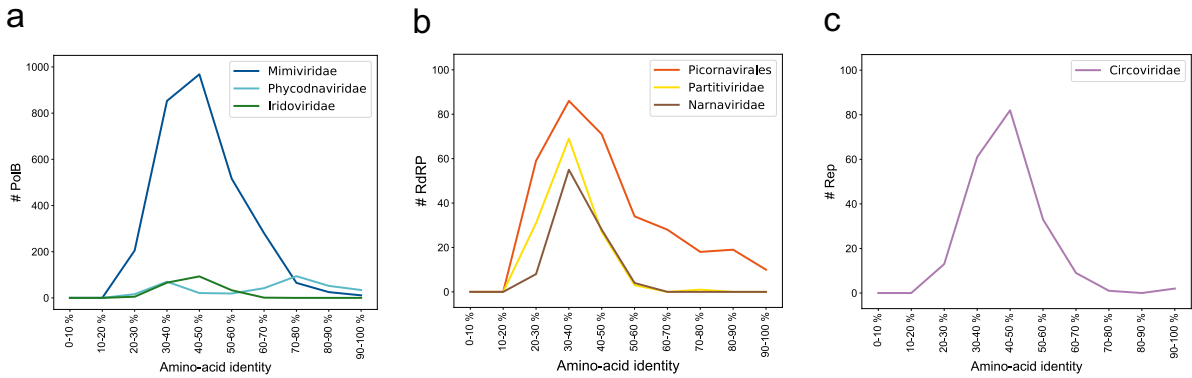
# Figure 4



## 1 Supplemental Information

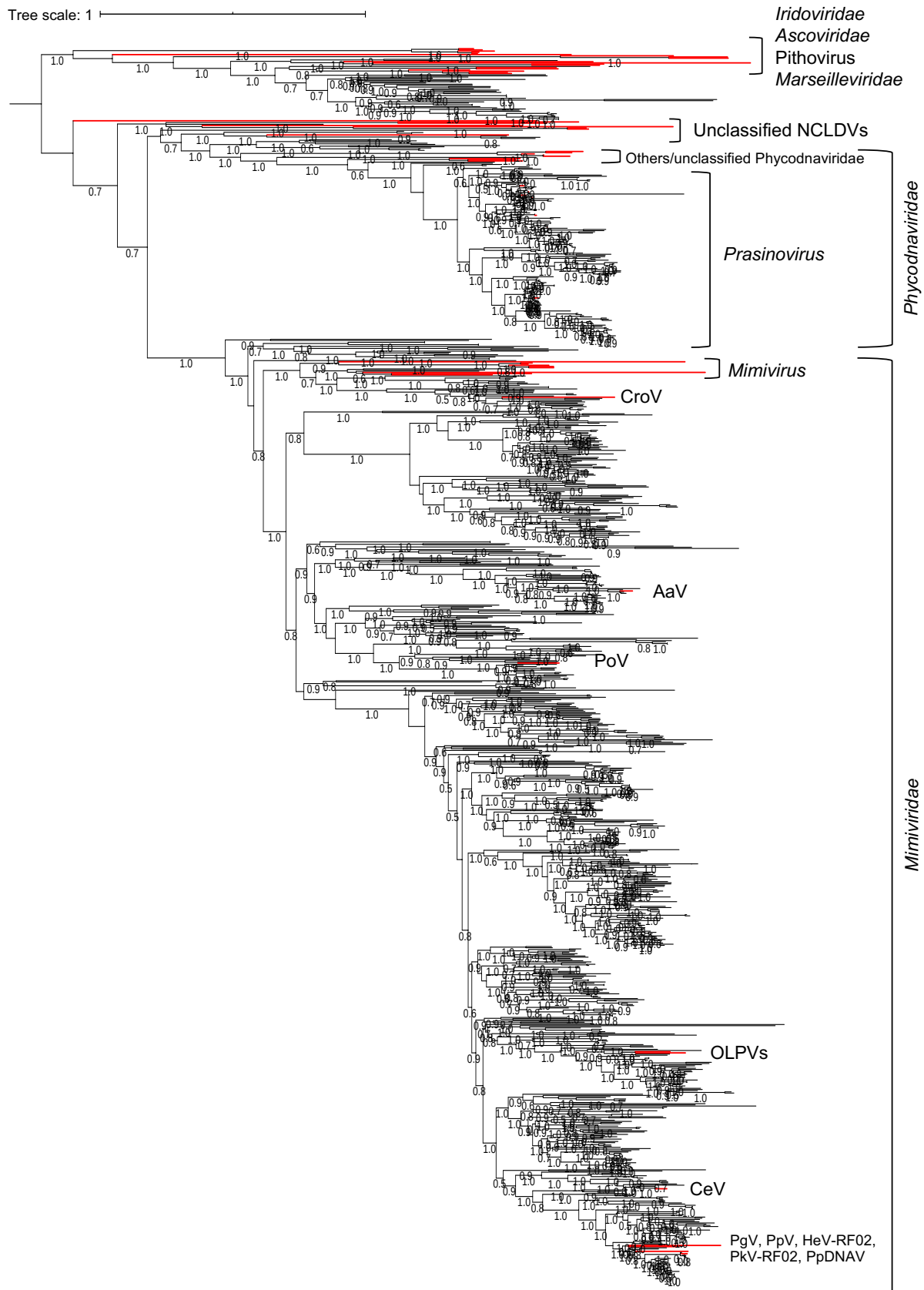
## 2 Supplemental Figures

3



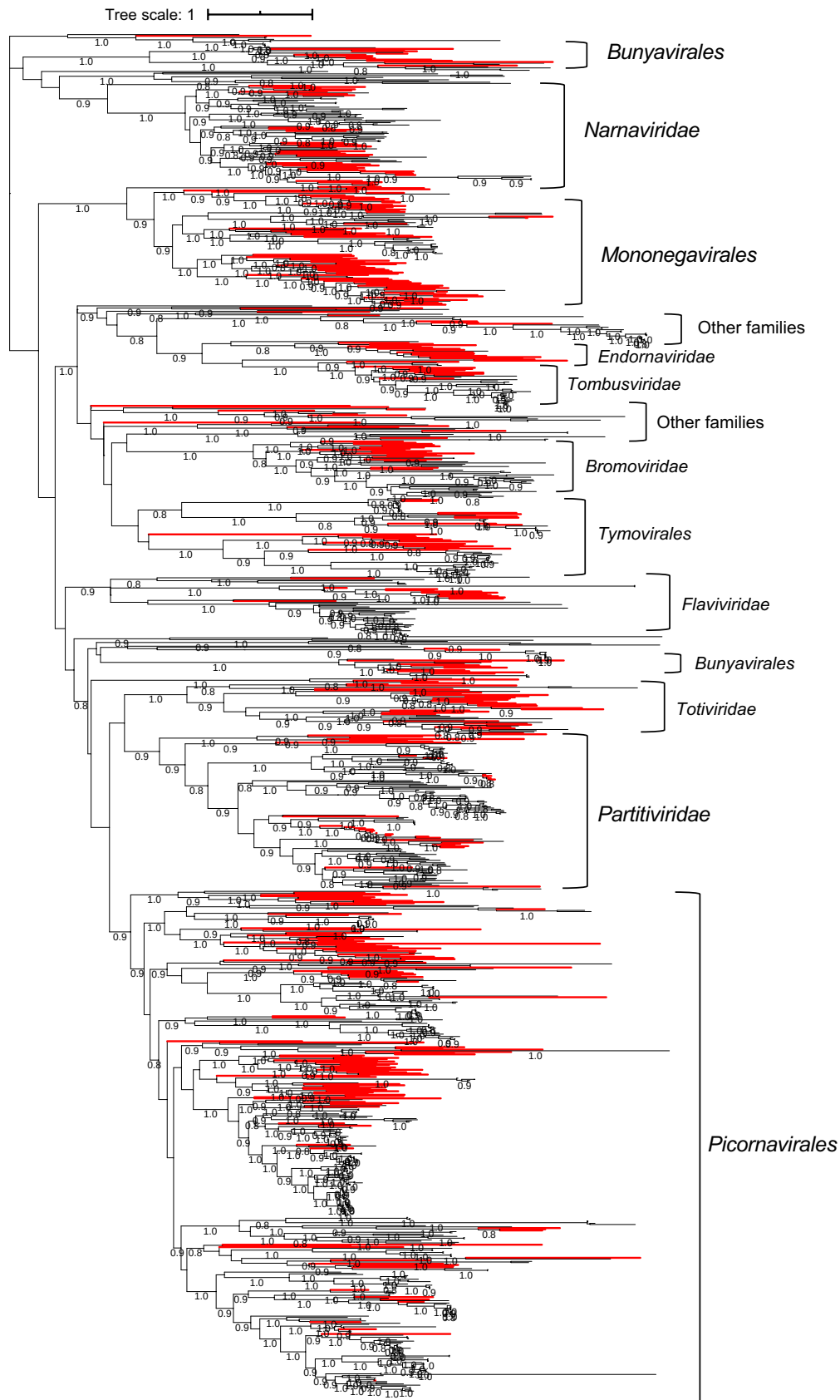
4

5 **Figure S1:** Distribution of the degree (%) of amino acid identity between environmental  
6 sequences and their best BLAST hits to reference sequences for nucleocytoplasmic large  
7 DNA viruses (NCLDV) (a), RNA viruses (b), and ssDNA viruses (c).

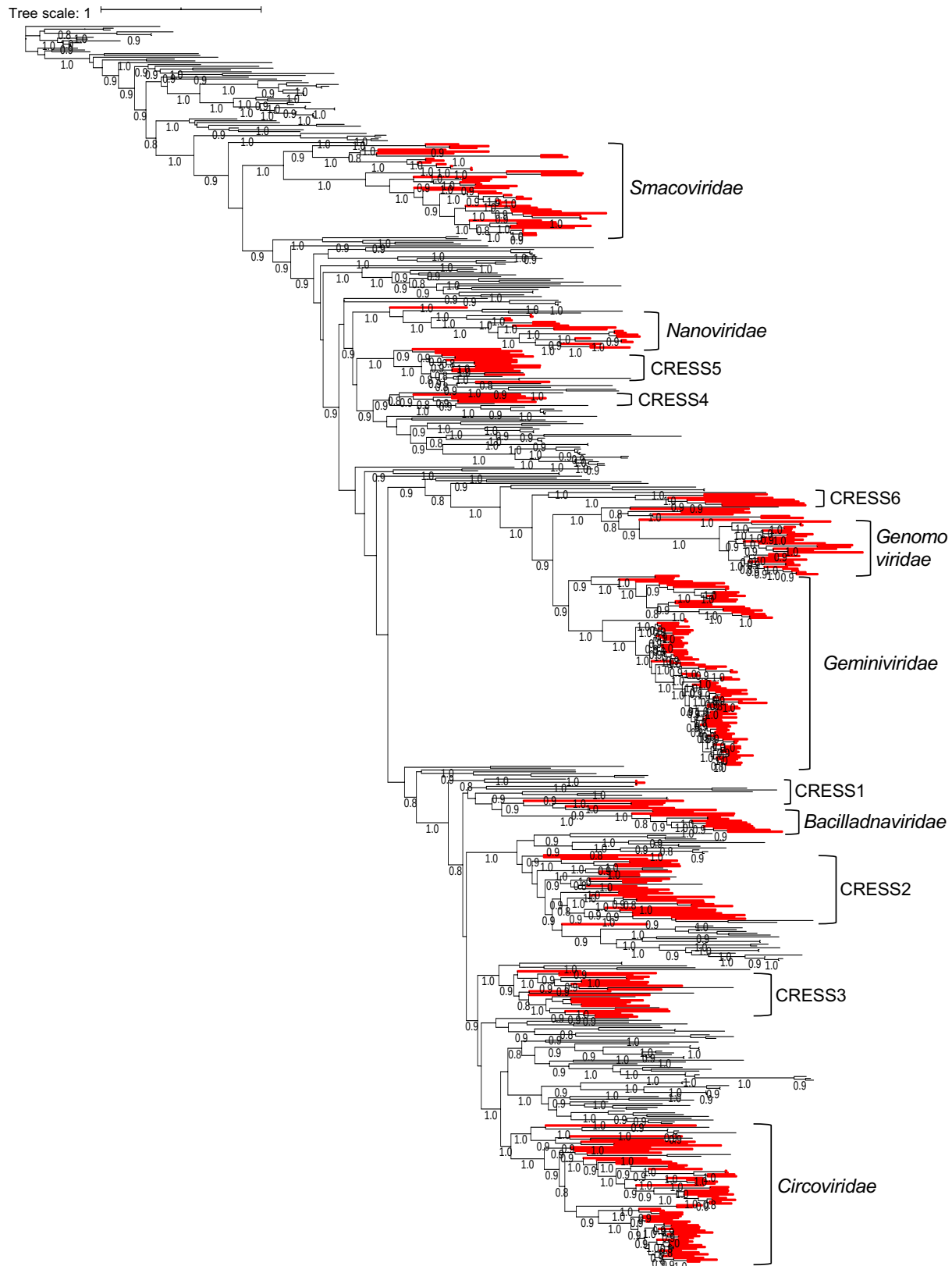


8  
9 **Figure S2:** Maximum likelihood phylogenetic trees for NCLDV DNA polymerase family B.  
10 Environmental sequences are shown in black and references in red. Approximate Shimodaira–  
11 Hasegawa (SH)-like local support values greater than 0.8 are shown. Scale bar indicates one  
12 change per site.

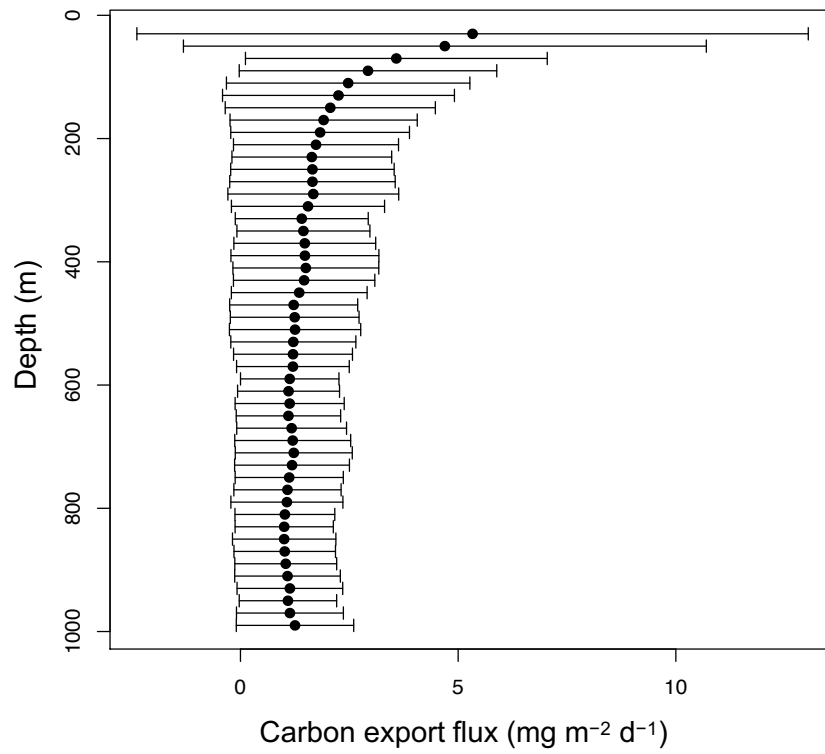




13  
14 **Figure S3:** Unrooted maximum likelihood phylogenetic trees for RNA virus RNA-dependent  
15 RNA polymerase. Environmental sequences are shown in black and references in red.  
16 Approximate SH-like local support values greater than 0.8 are shown. Scale bar indicates one  
17 change per site.

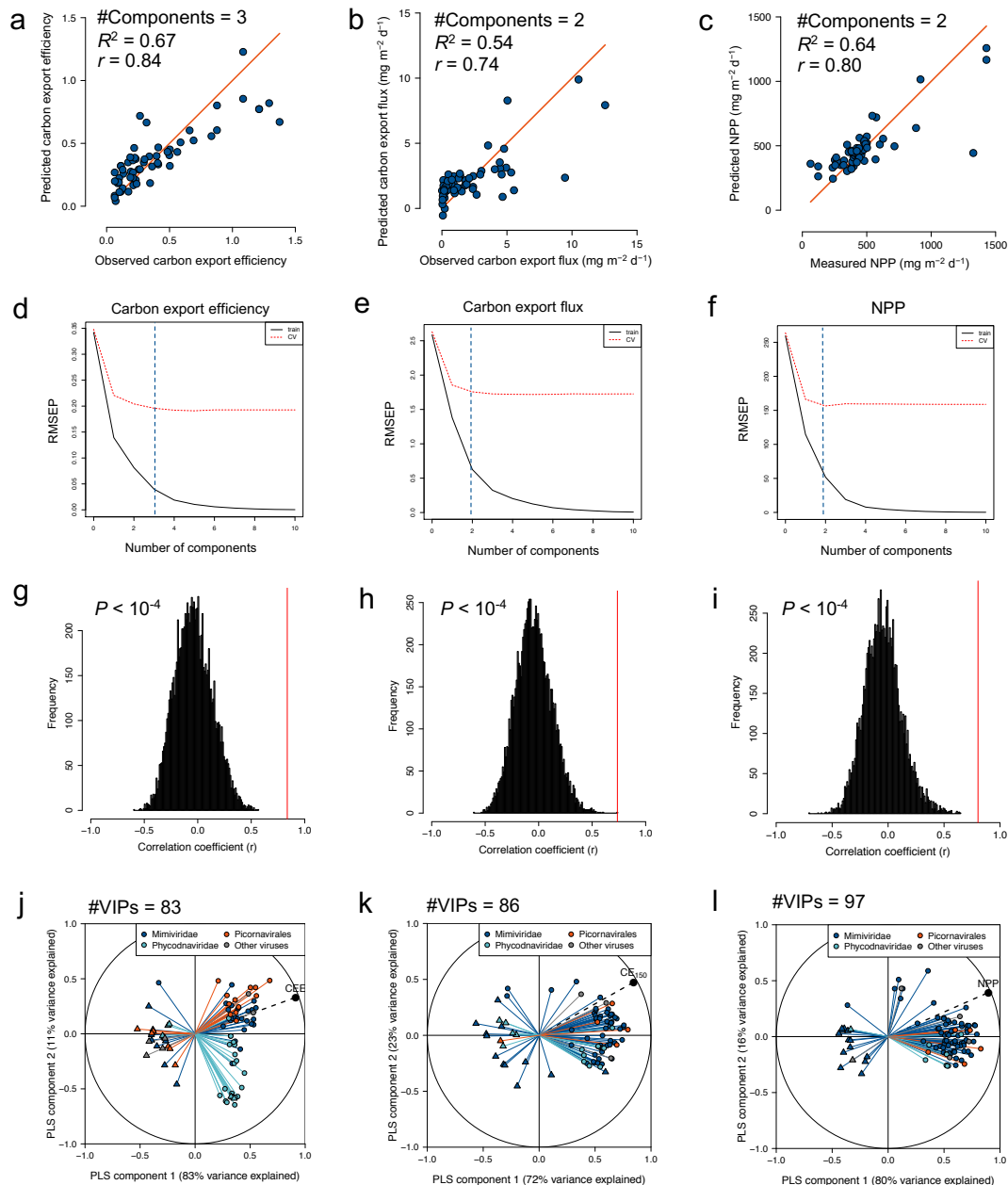


18  
19 **Figure S4:** Unrooted maximum likelihood phylogenetic trees for ssDNA virus replication-  
20 associated protein. Environmental sequences are shown in black and references in red.  
21 Approximate SH-like local support values greater than 0.8 are shown. Scale bar indicates one  
22 change per site.  
23

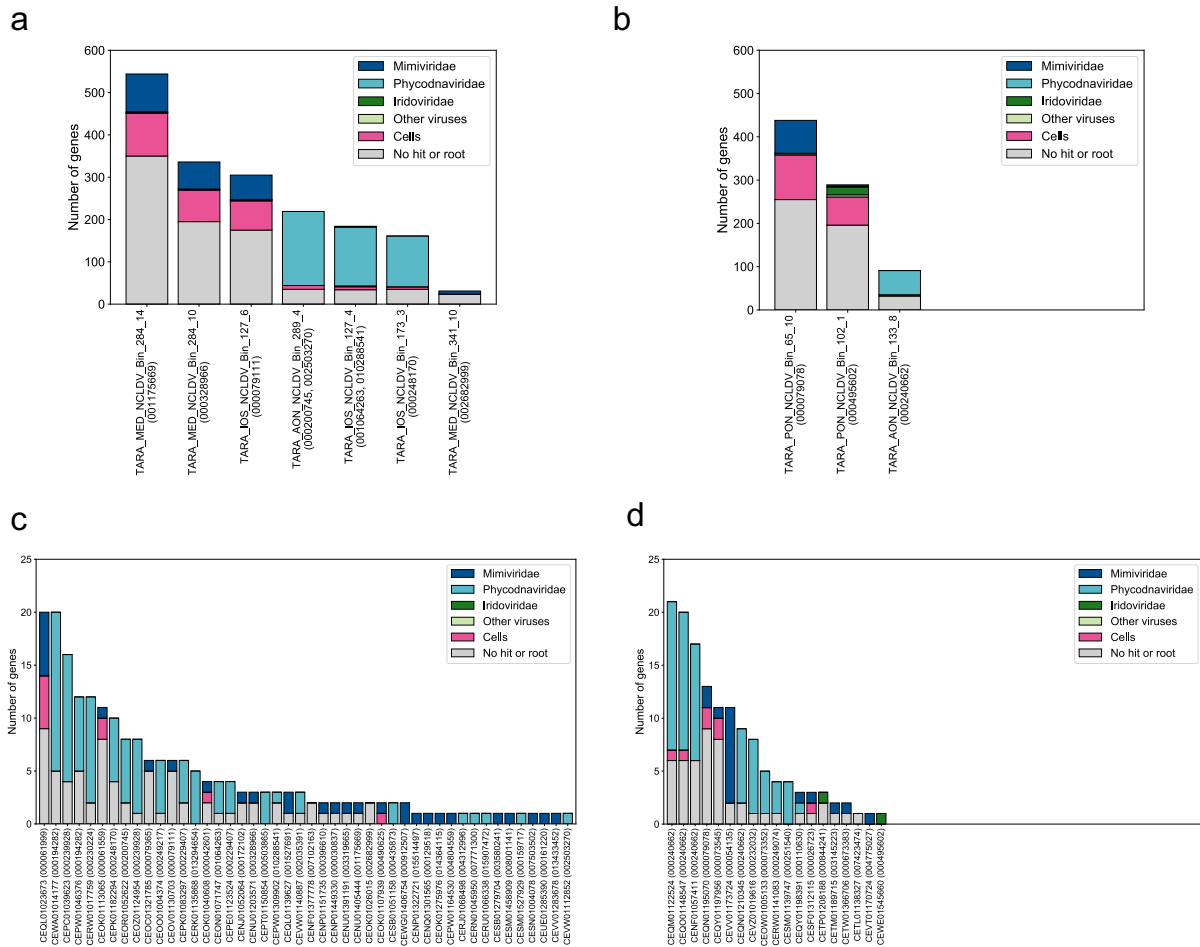


24  
25  
26  
27

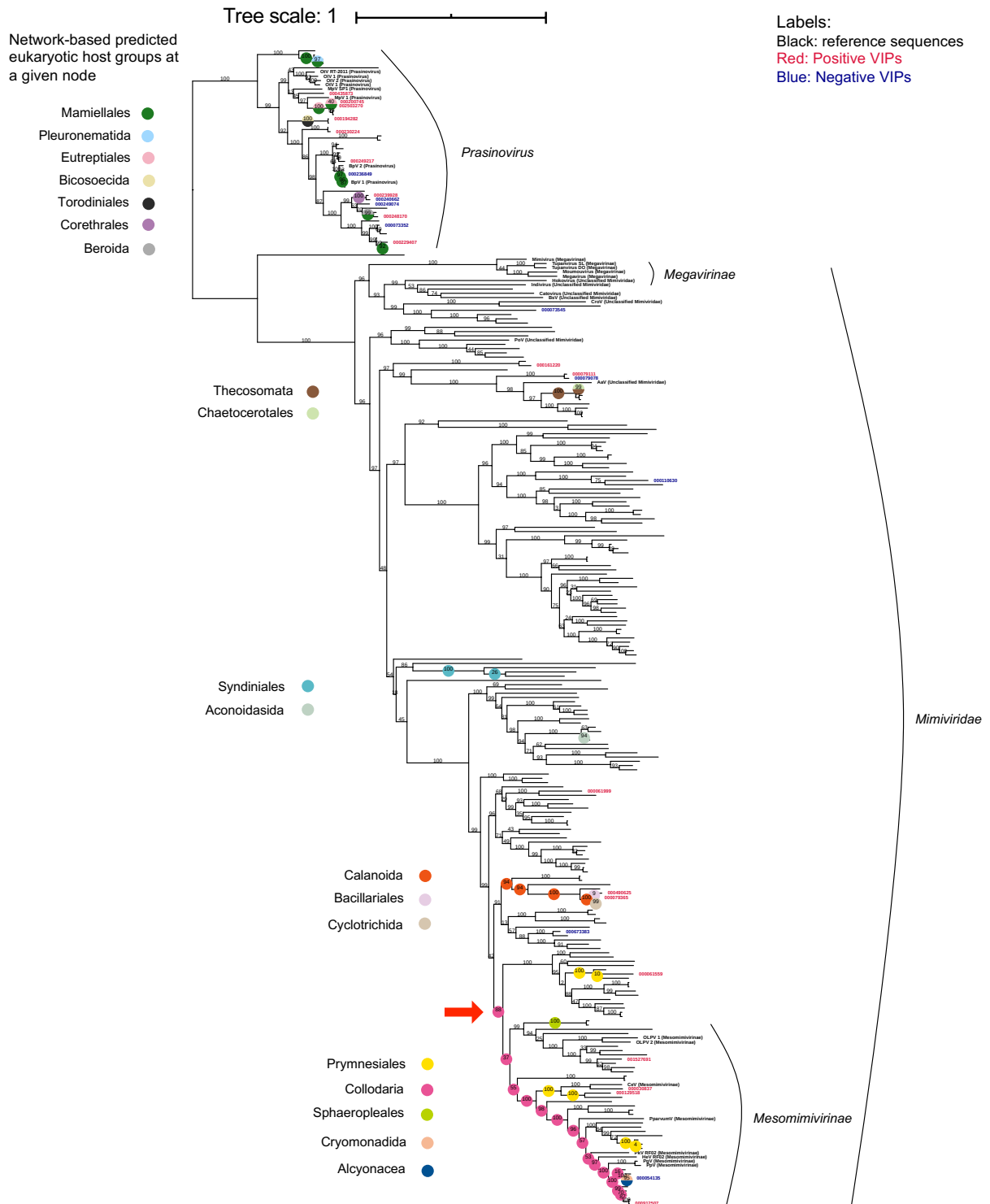
**Figure S5:** Variation in carbon export flux ( $\text{mg m}^{-2} \text{d}^{-1}$ ) across sampling depths in the water column. Dots are average values, and horizontal lines represent standard deviations.



28  
 29 **Figure S6:** The results of PLS regressions using relative abundance profiles of viral marker-  
 30 genes to explain the variance of carbon export efficiency (CEE) (**a, d, g, j**), carbon export flux  
 31 at 150 meters (CE<sub>150</sub>) (**b, e, h, k**), and net primary production (NPP) (**c, f, i, l**). **a–c** Bivariate  
 32 plots between predicted and observed response values in a leave-one-out cross-validation test.  
 33 The red diagonal line shows the theoretical curve for perfect prediction. **d–f** Variation in root  
 34 mean squared error of predictions (RMSEP) for the training set (solid black line) and cross-  
 35 validation set (red dashed line) across the number of components. Blue dashed line shows the  
 36 number of components selected for the analysis. **g–i** Results of the permutation tests ( $n =$   
 37 10,000) supporting the significance of the association between viruses and the response  
 38 variable. The histograms show the distribution of Pearson correlation coefficients obtained  
 39 from PLS models reconstructed based on the permuted response variable and red line show  
 40 the non-permuted response variable. **j–l** Pearson correlation coefficients between the  
 41 response variable and abundance profiles of viruses with VIP scores  $> 2$  (VIPs) with the first  
 42 two components in the PLS regression model using all samples. Viruses with positive  
 43 regression coefficients are shown with circles, and those with negative coefficients are shown  
 44 with triangles.

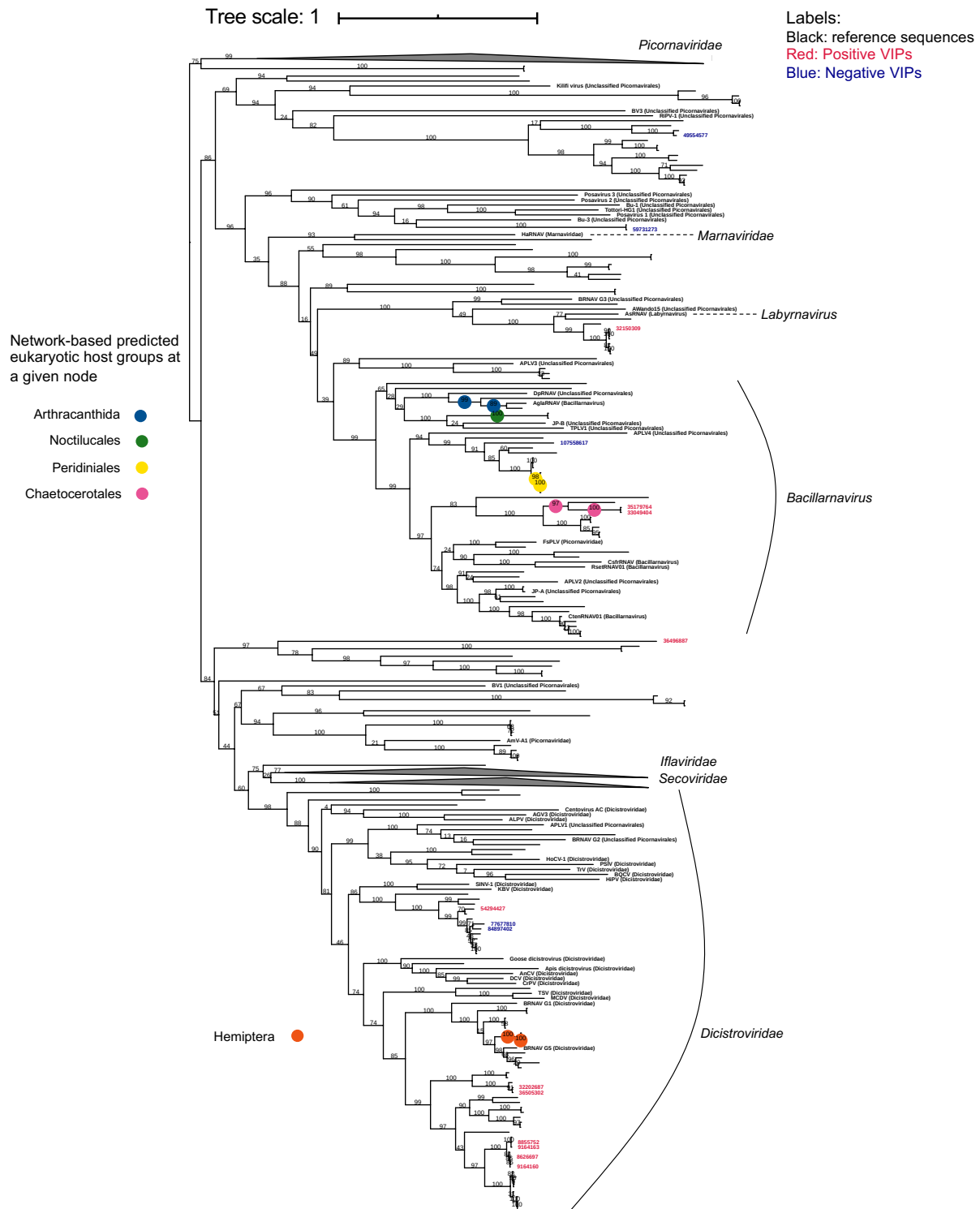


45  
 46 **Figure S7:** Taxonomic composition of genes predicted in viral genome fragments encoding  
 47 NCLDV PolBs positively (a and c) or negatively (b and d) associated with CEE (VIP score >  
 48 2). **a and b** Metagenome-assembled genomes (MAGs) derived from samples filtered to retain  
 49 particles of sizes > 0.8  $\mu\text{m}$ . **c and d** Contigs derived from samples filtered to retain particles  
 50 between 0.2  $\mu\text{m}$  and 3  $\mu\text{m}$  in size. Taxonomic annotations were performed as described in  
 51 Methods.

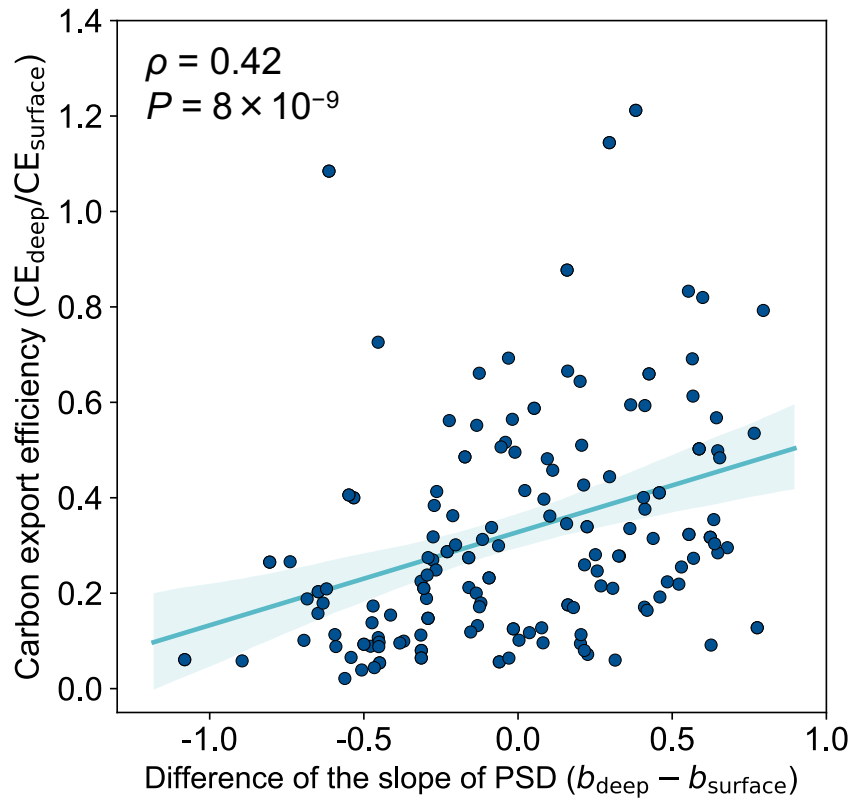


52  
53  
54  
55  
56  
57  
58  
59  
60

**Figure S8:** Phylogenetic positions of NCLDV PolBs associated with CEE and network-based predicted eukaryotic host groups. The unrooted maximum likelihood phylogenetic tree contains environmental (labeled in red if VIP score > 2 and the regression coefficient is positive, labeled in blue if negative) and reference (labeled in black) sequences of *Prasinovirus* and *Mimiviridae* PolBs. The approximate SH-like local support values are shown in percentages at nodes, and the scale bar indicates one change per site. Host groups predicted at nodes are shown with colored circles. The red arrow points to a clade of viruses predicted to infect Prymnesiales.



61  
 62 **Figure S9:** Phylogenetic position of *Piconavirales* RdRPs associated with CEE and network-  
 63 based predicted eukaryotic host groups. The unrooted maximum likelihood phylogenetic tree  
 64 contains environmental (labeled in red if VIP score > 2 and the regression coefficient is  
 65 positive, labeled in blue if negative) and reference (labeled in black) sequences of  
 66 *Piconavirales* RdRPs. The approximate SH-like local support values are shown in  
 67 percentages at nodes, and the scale bar indicates one change per site. Host groups predicted at  
 68 nodes are shown with colored circles.



69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80

**Figure S10:** Carbon export efficiency (CEE) is correlated with the change in the slope of particle size distribution (PSD) that occurred from the surface to deep (below the euphotic zone). Observed PSDs were fitted in the form  $n = ad^b$ , where  $n$  is the frequency of particles of a given size,  $d$  is the particle diameter, and  $a$  and  $b$  are parameters (as described by (Guidi et al., 2008)).  $b$ , the PSD slope, is a proxy for particles size. For example  $b = -5$  indicates presence of a large proportion of smaller particles, whereas  $b = -3$  indicates a preponderance of larger particles. A higher  $b$  value at deep compared to surface is suggestive of aggregation or presence of larger organisms at deep compare to surface. The blue line shows the regression line between CEE and the PSD slope difference between surface and deep. The shade around the regression line shows the 95% confidence interval.



81

82 **Supplemental Tables**

83

84 **Table S1: Viral lineages associated with CEE**

Viruses		VIPs	Positives VIPs	Negative VIPs
NCLDVs	Mimiviridae	34	25	9
	Phycodnaviridae	24	18	6
	Iridoviridae	2	0	2
	Other NCLDVs *	0	0	0
	Total	60	43	17
RNA viruses	Picornavirales (ssRNA+)	19	13	6
	Partitiviridae (dsRNA)	1	1	0
	Narnaviridae (ssRNA+)	0	0	0
	Other families	2*	0	2
	Unclassified	0	0	0
	RNA viruses	0	0	0
	Total	22	14	8
ssDNA viruses	Circoviridae	1	1	0
	Geminiviridae	0	0	0
	Nanoviridae	0	0	0
	Unclassified	0	0	0
	ssDNA viruses	0	0	0
	Total	1	1	0
All		83	58	25

85

86

\* Two Hepeviridae (ssRNA+).

87 **Table S2: Assembly statistics for NCLDV metagenome-assembled genomes and**  
 88 **corresponding VIPs**

Metagenome-assembled genome	#contigs	N50	L50	Min	Max	Sum	VIPs OTUs (OM-RGC.v1 ID)
TARA_IOS_NCLDV_Bin_127_6	14	21,642	5	8,581	35,822	267,607	PolB 000079111
TARA_IOS_NCLDV_Bin_173_3	12	12,913	3	2,807	34,517	108,412	PolB 000248170
TARA_MED_NCLDV_Bin_284_10	34	10,936	10	2,580	29,722	298,760	PolB 000328966
TARA_MED_NCLDV_Bin_284_14	43	14,837	11	2,756	27,607	439,843	PolB 001175669
TARA_IOS_NCLDV_Bin_127_4	26	5,734	10	2,560	8,505	133,765	PolB 001064263 and 010288541
TARA_AON_NCLDV_Bin_289_4	17	9,468	5	3,044	26,201	153,728	PolB 000200745 and 002503270
TARA_MED_NCLDV_Bin_341_10	5	7,800	2	2,534	7,941	30,478	PolB 002682999
TARA_PON_NCLDV_Bin_65_10	35	13,866	11	3,781	43,080	382,455	PolB 000079078
TARA_PON_NCLDV_Bin_102_1	53	4,608	18	2,606	11,485	239,832	PolB 000495602
TARA_AON_NCLDV_Bin_133_8	8	7,204	3	2,686	10,349	51,009	PolB 000240662

89 N50: The length of the contigs for which half of the assembly size is contained in contigs with a length greater than N50.  
 90 L50: Number of contigs (or scaffolds) with a size greater or equal to N50.  
 91

92 **Table S3: Host predictions per viral and host group for 83 VIPs based on**  
93 **phylogeny, co-occurrence analysis, and genomic context**

Virus-host relationship	Positive VIPs	Negative VIPs	Total
NCLDV-Mamiellales	10	4	14
NCLDV-Prymnesiales	5	1	6
NCLDV-Pelagophyceae	2	1	3
NCLDV-No prediction	26	11	37
RNA virus-Copepoda	7	2	9
RNA virus-Chaetocerotales	2	0	2
RNA virus-Labyrinthulomycetes	1	0	1
RNA virus-No prediction	4	6	10
ssDNA virus-Copepoda	1	0	1
Total	58	25	83

94



105 **Table S5: Statistics for the FlashWeave co-occurrence graphs**

Viral marker gene	Planktonic size fraction*	#Samples	#Viral OTUs	#Eukaryotic OTUs	#Edges in graph	#Virus-to-eukaryote edges	#Viruses connected to a eukaryote (%)
NCLDVs PolB	Piconano	99	2269	4936	20934	3594	1735 (76)
	Nano	51	1775	1872	6704	1027	721 (41)
	Micro	92	2205	2524	12189	2101	1299 (59)
	Meso	95	2238	2250	11624	1796	1126 (50)
RNA viruses RdRP	Piconano	60	125	4484	10754	446	122 (98)
	Nano	36	53	1768	2659	124	46 (87)
	Micro	62	124	2407	5351	367	117 (94)
	Meso	62	48	2100	4329	116	42 (88)
ssDNA viruses Rep	Piconano	60	64	4484	10577	205	63 (98%)
	Nano	36	1	1768	2563	2	1 (100%)
	Micro	62	4	2407	5086	9	4 (100%)
	Meso	62	8	2100	4242	24	8 (100%)

\* Pico: 0.8 to 5  $\mu\text{m}$ , Nano: 5 to 20  $\mu\text{m}$ , Micro: 20 to 180  $\mu\text{m}$ , Meso: 180 to 2000  $\mu\text{m}$

106  
107

108 **Table S6: Functional differences between eukaryotes found to be best connected to**  
 109 **VIPs and non-VIPs**

Functional trait	Positive VIPs ( <i>n</i> = 50)		Non-VIPs ( <i>n</i> = 983)		<i>P</i> -value (Fisher's exact test, two sided)	Adjusted <i>P</i> - value (BH) ( <i>Q</i> )
	Presence	Absence	Presence	Absence		
Chloroplast	20	30	276	690	0.109	0.164
Silicification	11	39	60	920	0.000	0.001
Calcification	1	49	30	950	1.000	1.000
Functional trait	Negative VIPs ( <i>n</i> = 21)		Non-VIPs ( <i>n</i> = 983)		<i>P</i> -value (Fisher's exact test, two sided)	Adjusted <i>P</i> - value (BH) ( <i>Q</i> )
	Presence	Absence	Presence	Absence		
Chloroplast	3	17	276	690	0.218	0.655
Silicification	0	21	60	920	0.632	0.947
Calcification	0	21	30	950	1.000	1.000

110

111 **Table S7: Functional differences between eukaryotes found to be best connected to**  
112 **positive and negative VIPs**

Functional trait	Positive VIPs (n = 50)		Negative VIPs (n = 21)		<i>P</i> -value (Fisher exact test, two sided)	Adjusted <i>P</i> - value (BH) ( <i>Q</i> )
	Presence	Absence	Presence	Absence		
Chloroplast	20	30	3	17	0.053	0.079
Silicification	11	39	0	21	0.027	0.079
Calcification	1	49	0	21	1.000	1.000

113

## 114 **Transparent Methods**

### 115 **Data context**

116 We used publicly available data generated in the framework of the *Tara* Oceans expedition.  
117 Single-copy marker-gene sequences for NCLDVs and RNA viruses were identified from two  
118 gene catalogs: the Ocean Microbial Reference Gene Catalog (OM-RGC) and the Marine Atlas  
119 of *Tara* Oceans Unigenes (MATOU). The viral marker-gene read count profiles used in our  
120 study are as previously reported for prokaryotic-sized metagenomes (size fraction 0.2–3  $\mu\text{m}$ )  
121 (Sunagawa et al., 2015) and eukaryotic-sized metatranscriptomes (Carradec et al., 2018).  
122 Eukaryotic plankton samples (the same samples were used for metatranscriptomes,  
123 metagenomes and 18S rRNA V9-meta-barcodes) were filtered for categorization into the  
124 following size classes: piconano (0.8–5  $\mu\text{m}$ ), nano (5–20  $\mu\text{m}$ ), micro (20–180  $\mu\text{m}$ ), and meso  
125 (180–2,000  $\mu\text{m}$ ). For eukaryotic 18S rRNA V9 OTUs (de Vargas et al., 2015), we used an  
126 updated version of the data that included functional trait annotations (chloroplast-bearing,  
127 silicified, and calcified organisms) of V9 OTUs. Occurrence profiles are compositional  
128 matrices in which gene occurrence are expressed as unnormalized (V9 meta-barcode data) or  
129 gene-length normalized (shotgun data) read counts. Indirect measurements of carbon export  
130 ( $\text{mg m}^{-2} \text{d}^{-1}$ ) in 5-m increments from the surface to a 1,000-m depth were taken from Guidi et  
131 al. (Guidi et al., 2016) The original measurements were derived from the distribution of  
132 particle sizes and abundances collected using an underwater vision profiler. These raw data  
133 are available from PANGEA (Picheral et al., 2014). Net primary production (NPP) data were  
134 extracted and averaged from 8-day composites of the vertically generalized production model  
135 (VGPM) (Behrenfeld and Falkowski, 1997) for the week of sampling. Thus, in this study, the  
136 comparisons between NPP and other parameters were not made at the same time point. This



137 might have affected the results of the regression analysis, especially if there were any short-  
138 term massive bloom events, although there was no bloom signal during the sampling period.

### 139 **Carbon export, carbon export efficiency, and particle size distribution**

140 Carbon flux profiles ( $\text{mg m}^{-2} \text{d}^{-1}$ ) were estimated based on particle size distributions and  
141 abundances. The method used for carbon flux estimation was previously calibrated comparing  
142 sediment trap measurement and data from imaging instruments (Guidi et al., 2008). Carbon  
143 flux values from depths of 30 to 970 meters were divided into 20-m bins, each obtained by  
144 averaging the carbon flux values from the designated 20 m in profiles gathered during  
145 biological sampling within a 25-km radius over 24 h when less than 50% of data were missing  
146 (Figure S5). Carbon export (CE) was defined as the carbon flux at 150 m (Guidi et al., 2016).  
147 Carbon export efficiency was calculated as follows:  $\text{CEE} = \text{CE}_{\text{deep}}/\text{CE}_{\text{surface}}$  (Buesseler and  
148 Boyd, 2009). To compare stations with different water column structures, we defined  $\text{CE}_{\text{surface}}$   
149 as the maximum CE (in a 20 m window) within the first 150 m.  $\text{CE}_{\text{deep}}$  is the average CE (also  
150 in a 20 m window) 200 m below this maximum. The 150 m limit serves as a reference point  
151 to automatize the calculation of  $\text{CE}_{\text{surface}}$  and  $\text{CE}_{\text{deep}}$ . The 150m-depth layer was selected  
152 because often used as a reference depth for drifting sediment trap and because most of the  
153 deep chlorophyll maximum (DCM) were shallower except at two (stations 98 (175 m) and  
154 100 (180 m)). The maximum  $\text{CE}_{\text{surface}}$  for these two stations was above 150 m. The sampling  
155 strategy of *Tara* Oceans was designed to study a variety of marine ecosystems and to target  
156 well-defined meso- to large-scale features (based on remote-sensing data). Therefore, this  
157 strategy avoided sampling water with important lateral inputs. Nevertheless, the possibility of  
158 having locations with potential lateral transport cannot be excluded.

159 We obtained the particle size distribution (PSD) profiles generated by the *Tara* Oceans  
160 expedition and computed the PSD slope at each depth for all profiles. The slope value  
161 (denoted “*b*”) is used as the descriptor of the particle size distribution as defined in a previous

162 work (Guidi et al., 2009). For example,  $b = -5$  indicates the presence of a large proportion of  
163 smaller particles, whereas  $b = -3$  indicates a preponderance of larger particles. We averaged  
164 the slope values at each sampling site in the same way as for carbon export flux.

## 165 **Identification of viral marker genes from ocean gene catalogs**

166 Viral genes were collected from two gene catalogs: OM-RGC version 1 and MATOU.  
167 Sequences in these two gene catalogs are representatives of clusters of environmental  
168 sequences (clustered at 95% nucleotide identity). The OM-RGC data were taxonomically re-  
169 annotated, with the NCBI reference tree used to determine the last common ancestor modified  
170 to reflect the current classification of NCLDV (Carradec et al., 2018). We automatically  
171 classified viral gene sequences as eukaryotic or prokaryotic according to their best BLAST  
172 score against viral sequences in the Virus-Host Database (Mihara et al., 2016). DNA  
173 polymerase B (PolB), RNA-dependent RNA polymerase (RdRP), and replication-associated  
174 protein (Rep) genes were used as markers for NCLDVs, RNA viruses, and ssDNA viruses,  
175 respectively. For PolB, reference proteins from the NCLDV orthologous gene cluster  
176 NCV0G0038 (Yutin et al., 2009) were aligned using MAFFT-*linsi* (Katoh and Standley,  
177 2013). A hidden Markov model (HMM) profile was constructed from the resulting alignment  
178 using *hmmbuild* (Eddy, 2011). This PolB HMM profile was searched against OM-RGC amino  
179 acid sequences and translated MATOU sequences annotated as NCLDVs, and sequences  
180 longer than 200 amino acids that had hits with  $E$ -values  $< 1 \times 10^{-5}$  were selected as putative  
181 PolBs. RdRP sequences were chosen from the MATOU catalog as follows: sequences  
182 assigned to Pfam profiles PF00680, PF00946, PF00972, PF00978, PF00998, PF02123,  
183 PF04196, PF04197, or PF05919 and annotated as RNA viruses were retained as RdRPs. For  
184 Rep, we reconstructed an HMM profile using a comprehensive set of reference sequences  
185 (Kazlauskas et al., 2018) and searched this profile against the translated MATOU sequences

186 annotated as ssDNA viruses. We kept sequences that had hits with  $E$ -values  $< 1 \times 10^{-5}$  and  
187 removed those that contained frameshifts.

188 The procedure above identified 3,486 PolB sequences in the metagenomic samples and  
189 respectively 975, 388, and 299 RdRP, PolB, and Rep sequences in the metranscriptomes.

## 190 **Testing for associations between viruses with CEE, CE<sub>150</sub>, and NPP**

191 To test for associations between occurrence of viral marker genes and CEE, CE<sub>150</sub>, and NPP,  
192 we proceeded as follows. Samples with CEE values greater than one and with  $Z$ -score greater  
193 than two were considered as outliers and removed (this removed the two samples from station  
194 68). Only marker genes represented by at least two reads in five or more samples were  
195 retained (lowering this minimal number of required samples down to three or four did not  
196 improve the PLS regression model). To cope with the sparsity and composition of the data,  
197 gene-length normalized read count matrices were center log-ratio transformed, separately for  
198 ssDNA viruses, RNA viruses and NCLDV. We next selected genes with Spearman  
199 correlation coefficients with CEE, CE<sub>150</sub> or NPP greater than 0.2 or smaller than  $-0.2$  (zero  
200 values were removed). To assess the association between these marker genes and CEE, we  
201 used partial least square (PLS) regression analysis. The number of components selected for  
202 the PLS model was chosen to minimize the root mean square error of prediction ([Figure S6](#)).  
203 We assessed the strength of the association between carbon export (the response variable) and  
204 viral marker genes occurrence (the explanatory variable) by correlating leave-one-out cross-  
205 validation predicted values with the measured carbon export values. We tested the  
206 significance of the correlation by comparing the original Pearson coefficients between  
207 explanatory and response variables with the distribution of Pearson coefficients obtained from  
208 PLS models reconstructed based on permuted data (10,000 iterations). We estimated the  
209 contribution of each gene (predictor) according to its variable importance in the projection  
210 (VIP) score derived from the PLS regression model using all samples. The VIP score of a

211 predictor estimates its contribution in the PLS regression. Predictors with high VIP scores (>  
212 2) were assumed to be important for the PLS prediction of the response variable.

### 213 **Phylogenetic analysis**

214 Environmental PolB sequences annotated as NCLDV s were searched against reference  
215 NCLDV PolB sequences using BLAST. Environmental sequences with hits to a reference  
216 sequence that had > 40% identity and an alignment length greater than 400 amino acids were  
217 kept and aligned with reference sequences using MAFFT-*linsi*. Environmental RdRP  
218 sequences annotated as were translated into six frames of amino acid sequences, and reference  
219 RNA viruses RdRP sequences collected from the Virus-Host Database were searched against  
220 the Conserved Domain Database (CDD) using rpsBLAST. The resulting alignment was used  
221 to trim reference and environmental RdRP sequences to the conserved part corresponding to  
222 the domain, CDD: 279070, before alignment with MAFFT-*linsi*. Rep sequences annotated as  
223 ssDNA viruses were treated similarly. PolB, RdRP, and Rep multiple sequence alignments  
224 were manually curated to discard poorly aligned sequences. Phylogenetic trees were  
225 reconstructed using the the *build* function of ETE3 (Huerta-Cepas et al., 2016) of the  
226 GenomeNet TREE tool (<https://www.genome.jp/tools-bin/ete>). Columns were automatically  
227 trimmed using *trimAl* (Capella-Gutiérrez et al., 2009), and trees were constructed using  
228 FastTree with default settings (Price et al., 2009).

229 A similar procedure was applied for the trees used in the hosts prediction analysis albeit  
230 selecting sequences for the Phycodnaviridae/Mimiviridae (Figure S8) and the Picornavirales  
231 (Figure S9) and removing the ones occurring in fewer than 10 samples, to reduce the size of  
232 the tree.

## 233 **Virus–eukaryote co-occurrence analysis**

234 We used FlashWeave (Tackmann et al., 2019) with Julia 1.2.0 to predict virus–host  
235 interactions based on their co-occurrence patterns. Read count matrices for the 3,486 PolBs,  
236 975 RdRPs, 299 Repls, and 18S rRNA V9 DNA barcodes obtained from samples collected at  
237 the same location were fed into FlashWeave. The 18S rRNA V9 data were filtered to retain  
238 OTUs with an informative taxonomic annotation. The 18S rRNA V9 OTUs and viral marker  
239 sequences were further filtered to conserve only those present in at least five samples.  
240 FlashWeave networks were learned for each of the four eukaryotic size fractions with the  
241 parameters ‘heterogenous’ = false and ‘sensitive’ = true, and edges receiving a weight > 0.2  
242 and a  $Q$ -value < 0.01 (the default) were retained. The number of samples per size fraction  
243 ranged between 51 and 99 for NCLDV and between 36 and 62 for RNA and ssDNA viruses.  
244 The number of retained OTUs per size fraction varied between 1,775 and 2,269 for NCLDVs  
245 and between 48 and 125 for RNA viruses (Table S5).

## 246 **Mapping of putative hosts onto viral phylogenies**

247 In our association networks, individual viral sequences were often associated with multiple  
248 18S rRNA V9 OTUs belonging to drastically different eukaryotic groups, a situation that can  
249 reflect interactions among multiple organisms but also noise associated with this type of  
250 analysis (Coenen and Weitz, 2018). To extract meaningful information from these networks,  
251 we reasoned as follows. We assumed that evolutionarily related viruses infect evolutionarily  
252 related organisms, similar to the case of phycodnaviruses (Clasen and Suttle, 2009). In the  
253 interaction networks, the number of connections between viruses in a given clade and the  
254 associated eukaryotic host group should accordingly be enriched compared with the number  
255 of connections with non-host organisms arising by chance. Following this reasoning, we  
256 assigned the most likely eukaryotic host group as follows. The tree constructed from viral  
257 marker-gene sequences (PolB, RdRP or Rep) was traversed from root to tips to visit every

258 node. We counted how many connections existed between leaves of each node and the V9-  
259 OTUs of a given eukaryotic group (order level). We then tested whether the node was  
260 enriched compared with the rest of the tree using Fischer's exact test and applied the  
261 Benjamini–Hochberg procedure to control the false discovery rate among comparisons of  
262 each eukaryotic taxon (order level). To avoid the appearance of significant associations driven  
263 by a few highly connected leaves, we required half of the leaves within a node to be  
264 connected to a given eukaryotic group. Significant enrichment of connections between a virus  
265 clade and a eukaryotic order was considered to be indicative of a possible virus–host  
266 relationship. We refer to the above approach, in which taxon interactions are mapped onto a  
267 phylogenetic tree of a target group using the organism's associations predicted from a species  
268 co-occurrence-based network, as TIM, for Taxon Interaction Mapper. This tool is available at  
269 <https://github.com/RomainBlancMathieu/TIM>. This approach can be extended to interactions  
270 other than virus–host relationships.

### 271 **Assembly of NCLDV metagenome-assembled genomes (MAGs)**

272 NCLDV metagenome-assembled genomes (MAGs) were assembled from *Tara* Oceans  
273 metagenomes corresponding to size fractions > 0.8  $\mu\text{m}$ . Metagenomes were first organized  
274 into 11 'metagenomic sets' based upon their geographic coordinates, and each set was co-  
275 assembled using MEGAHIT (Li et al., 2015) v.1.1.1. For each set, scaffolds longer than 2.5  
276 kbp were processed within the bioinformatics platform anvio (Eren et al., 2015) v.6.1  
277 following methodology described previously for genome-resolved metagenomics (Delmont et  
278 al., 2018). Briefly, we used the automatic binning algorithm CONCOCT (Alneberg et al.,  
279 2014) to identify large clusters of contigs using both sequence composition and differential  
280 coverage across metagenomes within the set. We then used HMMER (Eddy, 2011) v3.1b2 to  
281 search for a collection of eight NCLDV gene markers (Guglielmini et al., 2019), and  
282 identified NCLDV MAGs by manually binning CONCOCT clusters of interest using the

283 anvi'o interactive interface. The interface displayed hits for the eight gene markers alongside  
284 coverage values across metagenomes and GC-content. Finally, NCLDV MAGs were  
285 manually curated using the same interface, to minimize contamination as described previously  
286 (Delmont and Eren, 2016).

## 287 **Taxonomic composition of genes predicted in NCLDV genomes of VIPs**

288 VIP's PolB sequences were searched (using BLAST) against MAGs reconstructed from the  
289 metagenomes of the eukaryotic size fraction ( $> 0.8 \mu\text{m}$ ) and against contigs used to produce  
290 OM-RGCv1. Genome fragments covering 95% of the length of PolB VIPs with  $> 95\%$   
291 nucleotide identity were considered as originating from a same viral OTUs. Genes were  
292 predicted and annotated taxonomically with the same procedure described above  
293 (identification of viral marker genes). Genes contained in viral genome fragments and  
294 annotated as cellular organisms with amino acid identities  $> 60\%$  were manually inspected  
295 ([Supplemental Data 2](#)).

## 296 **Statistical test**

297 All the statistical significance assessments were performed with two-sided test.

## 298 **Supplemental References**

- 299 Alneberg, J., Bjarnason, B.S., Bruijn, I. de, Schirmer, M., Quick, J., Ijaz, U.Z., Lahti, L.,  
300 Loman, N.J., Andersson, A.F., and Quince, C. (2014). Binning metagenomic contigs by  
301 coverage and composition. *Nat. Methods* *11*, 1144–1146.
- 302 Behrenfeld, M.J., and Falkowski, P.G. (1997). Photosynthetic rates derived from satellite-  
303 based chlorophyll concentration. *Limnol. Oceanogr.* *42*, 1–20.
- 304 Buesseler, K.O., and Boyd, P.W. (2009). Shedding light on processes that control particle  
305 export and flux attenuation in the twilight zone of the open ocean. *Limnol. Oceanogr.* *54*,  
306 1210–1232.
- 307 Capella-Gutiérrez, S., Silla-Martínez, J.M., and Gabaldón, T. (2009). trimAl: a tool for  
308 automated alignment trimming in large-scale phylogenetic analyses. *Bioinforma. Oxf. Engl.*  
309 *25*, 1972–1973.

- 310 Carradec, Q., Pelletier, E., Silva, C.D., Alberti, A., Seeleuthner, Y., Blanc-Mathieu, R., Lima-  
311 Mendez, G., Rocha, F., Tirichine, L., Labadie, K., et al. (2018). A global ocean atlas of  
312 eukaryotic genes. *Nat. Commun.* *9*, 373.
- 313 Clasen, J.L., and Suttle, C.A. (2009). Identification of freshwater Phycodnaviridae and their  
314 potential phytoplankton hosts, using DNA pol sequence fragments and a genetic-distance  
315 analysis. *Appl. Environ. Microbiol.* *75*, 991–997.
- 316 Coenen, A.R., and Weitz, J.S. (2018). Limitations of Correlation-Based Inference in Complex  
317 Virus-Microbe Communities. *MSystems* *3*, e00084-18.
- 318 Delmont, T.O., and Eren, A.M. (2016). Identifying contamination with advanced visualization  
319 and analysis practices: metagenomic approaches for eukaryotic genome assemblies. *PeerJ* *4*,  
320 e1839.
- 321 Delmont, T.O., Quince, C., Shaiber, A., Esen, Ö.C., Lee, S.T., Rappé, M.S., McLellan, S.L.,  
322 Lückner, S., and Eren, A.M. (2018). Nitrogen-fixing populations of Planctomycetes and  
323 Proteobacteria are abundant in surface ocean metagenomes. *Nat. Microbiol.* *3*, 804–813.
- 324 Eddy, S.R. (2011). Accelerated Profile HMM Searches. *PLOS Comput. Biol.* *7*, e1002195.
- 325 Eren, A.M., Esen, Ö.C., Quince, C., Vineis, J.H., Morrison, H.G., Sogin, M.L., and Delmont,  
326 T.O. (2015). Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* *3*,  
327 e1319.
- 328 Guglielmini, J., Woo, A.C., Krupovic, M., Forterre, P., and Gaia, M. (2019). Diversification  
329 of giant and large eukaryotic dsDNA viruses predated the origin of modern eukaryotes. *Proc.*  
330 *Natl. Acad. Sci.* *116*, 19585–19592.
- 331 Guidi, L., Jackson, G.A., Stemmann, L., Miquel, J.C., Picheral, M., and Gorsky, G. (2008).  
332 Relationship between particle size distribution and flux in the mesopelagic zone. *Deep Sea*  
333 *Res. Part Oceanogr. Res. Pap.* *55*, 1364–1374.
- 334 Guidi, L., Stemmann, L., Jackson, G.A., Ibanez, F., Claustre, H., Legendre, L., Picheral, M.,  
335 and Gorsky, G. (2009). Effects of phytoplankton community on production, size, and export  
336 of large aggregates: A world-ocean analysis. *Limnol. Oceanogr.* *54*, 1951–1963.
- 337 Guidi, L., Chaffron, S., Bittner, L., Eveillard, D., Larhlimi, A., Roux, S., Darzi, Y., Audic, S.,  
338 Berline, L., Brum, J.R., et al. (2016). Plankton networks driving carbon export in the  
339 oligotrophic ocean. *Nature* *532*, 465.
- 340 Huerta-Cepas, J., Serra, F., and Bork, P. (2016). ETE 3: Reconstruction, analysis and  
341 visualization of phylogenomic data. *Mol. Biol. Evol.* msw046.
- 342 Katoh, K., and Standley, D.M. (2013). MAFFT multiple sequence alignment software version  
343 7: improvements in performance and usability. *Mol. Biol. Evol.* *30*, 772–780.
- 344 Kazlauskas, D., Varsani, A., and Krupovic, M. (2018). Pervasive Chimerism in the  
345 Replication-Associated Proteins of Uncultured Single-Stranded DNA Viruses. *Viruses* *10*.



- 346 Li, D., Liu, C.-M., Luo, R., Sadakane, K., and Lam, T.-W. (2015). MEGAHIT: an ultra-fast  
347 single-node solution for large and complex metagenomics assembly via succinct de Bruijn  
348 graph. *Bioinforma. Oxf. Engl.* *31*, 1674–1676.
- 349 Mihara, T., Nishimura, Y., Shimizu, Y., Nishiyama, H., Yoshikawa, G., Uehara, H., Hingamp,  
350 P., Goto, S., and Ogata, H. (2016). Linking Virus Genomes with Host Taxonomy. *Viruses* *8*,  
351 66.
- 352 Picheral, M., Searson, S., Taillandier, V., Bricaud, A., Boss, E., Stemmann, L., Gorsky, G.,  
353 Tara Oceans Consortium, C., and Tara Oceans Expedition, P. (2014). Vertical profiles of  
354 environmental parameters measured from physical, optical and imaging sensors during Tara  
355 Oceans expedition 2009-2013.
- 356 Price, M.N., Dehal, P.S., and Arkin, A.P. (2009). FastTree: Computing Large Minimum  
357 Evolution Trees with Profiles instead of a Distance Matrix. *Mol. Biol. Evol.* *26*, 1641–1650.
- 358 Sunagawa, S., Coelho, L.P., Chaffron, S., Kultima, J.R., Labadie, K., Salazar, G.,  
359 Djahanschiri, B., Zeller, G., Mende, D.R., Alberti, A., et al. (2015). Ocean plankton. Structure  
360 and function of the global ocean microbiome. *Science* *348*, 1261359.
- 361 Tackmann, J., Matias Rodrigues, J.F., and von Mering, C. (2019). Rapid Inference of Direct  
362 Interactions in Large-Scale Ecological Networks from Heterogeneous Microbial Sequencing  
363 Data. *Cell Syst.* *9*, 286-296.e8.
- 364 de Vargas, C., Audic, S., Henry, N., Decelle, J., Mahé, F., Logares, R., Lara, E., Berney, C.,  
365 Le Bescot, N., Probert, I., et al. (2015). Ocean plankton. Eukaryotic plankton diversity in the  
366 sunlit ocean. *Science* *348*, 1261605.
- 367 Yutin, N., Wolf, Y.I., Raoult, D., and Koonin, E.V. (2009). Eukaryotic large nucleo-  
368 cytoplasmic DNA viruses: Clusters of orthologous genes and reconstruction of viral genome  
369 evolution. *Virology* *6*, 223.
- 370



Click here to access/download  
**Supplemental File Sets**  
Supplemental\_Data\_1\_2.xlsx

