

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

**VIRUS-HOST INFECTION DYNAMICS OF MARINE SINGLE-CELLED EUKARYOTES RESOLVED FROM METATRANSCRIPTOMICS**

Mohammad Moniruzzaman<sup>a</sup>, Louie L. Wurch<sup>b</sup>, Harriet Alexander<sup>c</sup>, Sonya T. Dyhrman<sup>c</sup>,  
Christopher J. Gobler<sup>d</sup>, Steven W. Wilhelm<sup>a,1</sup>

<sup>a</sup> Department of Microbiology, The University of Tennessee, TN, 37996;

<sup>b</sup> Department of Biology, James Madison University, Harrisonburg, VA 22807;

<sup>c</sup> Department of Earth and Environmental Science and Lamont-Doherty Earth Observatory, Columbia University, NY 10964.

<sup>d</sup> School of Marine and Atmospheric Sciences, Stony Brook University, Stony Brook, NY 11794.

1. **Author for correspondence:** [wilhelm@utk.edu](mailto:wilhelm@utk.edu)

**Keywords:** Eukaryotic viruses, Brown tides, Virus diversity, picoeukaryotes, marine microbiology

27 **Abstract**

28 Metatranscriptomics has emerged as a tool in microbial ecology that can resolve the functional  
29 landscape of both prokaryotes and eukaryotes within a community. In this study, we extend the  
30 potential of metatranscriptomics to probe active virus infections and virus-host relationships in  
31 marine systems. Polyadenylation-selected RNA-seq data were examined from microbial  
32 communities in two productive marine environments: a brown tide bloom event dominated by  
33 *Aureococcus anophagefferens* in Quantuck Bay, NY, and a diatom-dominated plankton  
34 community in Narragansett Bay, RI. Active infections by diverse giant viruses (NCLDV) of  
35 algal and non-algal hosts were found at both sites. Ongoing infections of *A. anophagefferens* by  
36 a known *Mimiviridae* (AaV) were observed during both the peak and decline of the bloom.  
37 Bloom decline was also accompanied by increased activity for viruses other than AaV, including  
38 (+) ssRNA viruses. In Narragansett Bay, increased temporal resolution revealed active NCLDV  
39 with both ‘boom-and-bust’ as well as ‘steady-state infection’-like ecologies. Statistical co-  
40 occurrence examinations of the dsDNA, ssRNA and dsRNA markers within the data revealed a  
41 broad spectrum of statistically strong and significant virus-host relationships that included both  
42 known as well as novel interactions. Our approach offers a method for screening the diversity  
43 and dynamics of active viral infections in natural systems and develops links between viruses  
44 and their potential hosts *in situ*.

45

46 **Significance**

47 Viruses are important partners in ecosystem scale ecology, yet their study to date is primarily  
48 limited to single virus-host infection models in the laboratory or limited to “potential-actions”  
49 derived from metagenomics analyses. Using metatranscriptomic sequences from polyadenylated-  
50 RNA selected samples, we have simultaneously captured information regarding eukaryotic  
51 diversity and active infection by viruses with dsDNA genomes, resulting in a statistical  
52 opportunity to predict “*who is infecting whom*”. This approach further provides concurrent  
53 insight regarding viruses with ssRNA and dsRNA genomes, capturing dynamics for the  
54 communities of viruses infecting single-celled eukaryotes. Given the central role of these  
55 plankton in global scale processes, our efforts result in a transformational step-forward regarding  
56 the study of *in situ* virus-host interactions.

57

## 58 Introduction

59 Viruses that infect marine microbes are an integral component of aquatic ecosystems, with a  
60 diversity spectrum spanning the entire Baltimore classification scheme (1). The association of  
61 viruses with global-scale biogeochemistry, algal bloom termination events, and their impact on  
62 microbial community diversity has driven scientific research in virus ecology (2, 3). Amongst  
63 these predators, giant dsDNA viruses belonging to the Nucleocytoplasmic Large DNA Virus  
64 (NCLDV) group infect single-celled eukaryotes with diverse lifestyles (4) and are thought to be  
65 abundant in the world's oceans (5). Individually some of these viruses have been shown to be  
66 potential drivers of algal bloom collapse (6, 7). However, only a few NCLDV-host-systems with  
67 established ecological relevance have been identified. As a specific example the environmental  
68 hosts of the *Mimiviridae*, isolated using *Acanthamoeba* in the laboratory, are yet to be confirmed  
69 (8).

70 Along with the NCLDVs, RNA viruses also comprise a major fraction of the marine  
71 viroplankton, infecting organisms ranging from diatoms and dinoflagellates to fish (9). However,  
72 little is known about the ecology and host range of RNA viruses; the first RNA virus infecting a  
73 marine single-celled eukaryote was only described in 2003 (10). In addition, recent evidence  
74 suggests that a large number of novel ssDNA virus families possibly infect yet to be  
75 characterized marine phytoplankton and zooplankton (11). Collectively, these observations  
76 illustrate the strong need to develop *in situ* approaches that link the marine virosphere to their  
77 hosts within the microbial eukaryotes.

78 The marine ecosystem consists of complex interactions among diverse organisms and their  
79 viruses. While studying individual host-virus systems remain critical to understand the molecular  
80 basis of interactions, studying the overall contribution of viruses in a dynamic network of  
81 organisms is hindered by methodological limitations. Culture independent approaches to study  
82 viruses, and especially viral communities, are challenging: “viromes” – large metagenomic  
83 datasets enriched with viral sequences, are usually generated by size exclusion ( $\leq 0.22 \mu\text{m}$ ) of  
84 bacteria and small eukaryotes (2). This approach, however, largely removes the large virus  
85 particles that can range from 100 nM to  $\sim 1.5 \mu\text{M}$  (8). Moreover, by targeting DNA, these  
86 approaches examine only the presence of particles and not their activity. Additionally, RNA-  
87 containing virus particles must be targeted separately from DNA viruses, since common methods

88 for virus enumeration (using dsDNA intercalating stains) and DNA-based metagenomics  
89 approaches cannot detect them (12). Consequently, there is a need for new toolsets to  
90 complement the current approaches and yet overcome the aforementioned issues to provide a  
91 more comprehensive picture of the viral dynamics

92 Here we examined metatranscriptomes from two highly productive marine sites on the east  
93 coast of USA – Quantuck Bay, NY and Narragansett Bay, RI. Quantuck Bay experiences  
94 recurring ecosystem disruptive brown tide blooms caused by the pelagophyte *Aureococcus*  
95 *anophagefferens* (13) which are shaped by a giant virus (AaV) (14, 15). Narragansett Bay is a  
96 highly productive system with seasonal diatom blooms, but a poorly described eukaryotic virus  
97 community. By employing selection for polyadenylation prior to sequencing, we were able to  
98 focus on active virus infections within eukaryotes. By using time-series data, we were able to  
99 capture emergent relationships of putative virus-host pairings and their ecological dynamics.  
100 This approach also allowed us characterize viruses actively infecting eukaryotes with diverse  
101 nucleic acid genomes (ssDNA, ssRNA and dsRNA).

102

## 103 **Results and discussion**

### 104 **Temporal dynamics of active giant virus infections**

105 To identify NCLDV, we screened contig libraries generated at each study site for ten  
106 conserved NCLDV core genes (16). Reads from individual samples were mapped to the core  
107 gene contigs followed by library size normalization. At both sites, numerous contigs originating  
108 from NCLDV-specific Major Capsid Protein (MCP) were identified (Fig. 1). The abundance of  
109 reads mapped to MCP contigs was higher than the sum of mapped reads to all other NCLDV  
110 core gene contigs (Fig. 1) for all samples except QB-S3, confirming efforts suggesting MCP is a  
111 suitable marker for NCLDV diversity (Moniruzzaman et al. 2016) and that the MCP gene is  
112 highly expressed (17). Only distant homologs of MCP are present in *Poxviridae* (16) and there  
113 are no homologs in recently discovered Pandora- and Pithoviruses (8): to this end the ubiquity of  
114 this gene in all other NCLDV families makes it an excellent candidate for phylogenetic probing  
115 of metatranscriptomics data.

116 We placed the MCP contigs on a reference phylogenetic tree and studied their relative  
117 expression levels in terms of a metric that we defined as ‘rarefied counts per kilobase’ (RCK)  
118 (details in Materials and methods). Phylogenetic placement of the contigs demonstrated that  
119 NCLDV members from *Mimiviridae* and *Phycodnaviridae* were consistently present in both  
120 Quantuck Bay and Narragansett Bay. At both sites the highest number of contigs fell within the  
121 *Mimiviridae* family, followed by *Phycodnaviridae* (Fig. 2). A large number of contigs had strong  
122 phylogenetic affinity to AaV as well as other alga-infecting members of the *Mimiviridae* clade.  
123 Their presence and relative abundance in these field surveys demonstrate the *Mimiviridae* are an  
124 important component of the marine virosphere and are as active as the better-studied  
125 *Phycodnaviridae* group.

126 Brown tide bloom samples collected on June 14 (QB-S1) and June 16 (QB-S2) represented the  
127 bloom peak with an *Aureococcus* count of  $\sim 2.28 \times 10^6$  cells/mL and  $\sim 2.23 \times 10^6$  cells/mL,  
128 respectively. The third sample, collected on June 22, represented the early stage of bloom  
129 demise, with an *Aureococcus* count of  $\sim 1.91 \times 10^6$  cells/mL. We detected a persistent infection  
130 of *A. anophagefferens* by AaV across this sampling period. High stringency (similarity  $\geq 97\%$ )  
131 mapping of reads to the genome identified 1,368 and 604 reads that could be assigned to peak  
132 bloom samples QB-S1 and S2, respectively, after library size normalization, while 236 reads  
133 were mapped to the QB-S3 sample taken during bloom decline (Fig. S1). Across the entire  
134 genome, 15 AaV transcripts had more than 10 reads: roughly two thirds of these transcripts have  
135 no similarity to genes with currently known functions (ORFans) (Dataset S1). Highly expressed  
136 ORFans have also been recorded for Mimivirus: 17 of the top 20 most highly expressed genes  
137 were hypotheticals (17). These observations suggest these genes are active during infection by  
138 AaV and other NCLDVs, and represent important targets for future studies. Notably, the AaV  
139 MCP was among the most highly expressed functional genes, with 121 total reads mapped to this  
140 gene across the three *in situ* samples from Quantuck Bay. Both total reads mapped across the  
141 AaV genome (Fig. S1) and specifically to MCP gene (Fig. 2A) progressively declined  
142 throughout the sampling period, with the lowest number of reads mapped from S3. It may be that  
143 AaV activity was present, but reduced during the bloom decline stage - an observation that is  
144 supported by a recent study where AaV amplicons were only detected during the peak of the  
145 bloom (18). Overall these data reinforce the utility of MCP as a marker, since the MCP dynamics  
146 were consistent with data derived from the full AaV genomic analysis (Figure 2).

147 With five *in situ* samples over a period of approximately four weeks (Table S1), data from  
148 Narragansett Bay allowed us to observe the temporal dynamics of the NCLDV. Some members  
149 from *Phycodna*- and *Mimiviridae* clades showed persistent evidence of infection over a  
150 prolonged period, while ‘boom-bust’ like relationships (4) were possibly present for other  
151 members (Fig. 2B). For example, a number of MCP contigs were consistently expressed (within  
152 an order magnitude) between samples across time points (*e.g.*, blue arrows in Fig. 2B), an  
153 observation supporting the presence of infected hosts. While this scenario is consistent with a  
154 ‘slow-and-steady’ infection dynamic (19), it can also be explained by persistent infections of the  
155 plankton – where ongoing virus production does not necessarily lead to host (or at least total  
156 community) mortality (20). The expression of other phylogenetically distinct markers, however,  
157 reflected a ‘boom-and-bust’ like scenario (19), with the expression varying across several orders  
158 of magnitude between time points. One striking example of such putative ‘boom-and-bust’  
159 scenario was a contig in the non-algal *Mimiviridae* family, where expression decreased by two  
160 orders of magnitudes from May 16 to May 21 and May 30 (Fig. 2B, red arrow).

161

## 162 **Viruses infecting single-celled eukaryotes beyond the NCLDVs**

163 The marine virosphere is not limited to dsDNA viruses, as viruses containing all nucleic acid  
164 types (ss- and dsRNA as well as ssDNA) that infect marine single-celled eukaryotes have been  
165 described (9, 11). We extended our approach to detect the contigs that potentially originated  
166 from diverse RNA and DNA viruses other than NCLDVs. RNA viruses have a diverse size  
167 range, with Picornavirales particles as small as ~25-30 nm (21). Our sample collection method  
168 (Material and Methods) allowed detection of both ongoing virus infection (for DNA and RNA  
169 viruses) and cell-surface associated RNA viruses. It is important to mention that some (+)  
170 ssRNA viruses have poly-A tailed genomes (*e.g.*, Picorna- and Togaviruses) even outside the  
171 host (22). Therefore, owing to their nature, the (+) ssRNA viral diversity captured by this  
172 approach might reflect both actively replicating and surface bound viruses.

173 Within our analyses, 579 and 599 contigs from Quantuck Bay and Narragansett Bay,  
174 respectively, were assigned to viruses other than NCLDVs. The majority of these contigs  
175 originated from (+) ssRNA viruses, with the main contributors coming from a yet unclassified  
176 group of viruses in the Picornavirales order (23). Unclassified Picornavirales contigs represented

177 62% of the total non-NCLDV viral contigs for Quantuck Bay and 74% of this group for  
178 Narragansett Bay. Marine Picornavirales have been shown to infect diatoms (*e.g.*, *Chaetoceros*  
179 *sp.*, *Asterionellopsis glacialis* and *Rhizosolenia setigera*) (9) and a marine fungoid protist,  
180 *Aurantiotrychium* (24). The closest phylogenetic relative of this group is *Marnaviridae*, which  
181 currently have only one member – HaRNAV, that infects the marine raphidophyte *Heterosigma*  
182 *akashwo* (10). The second major group of (+) ssRNA viruses belonged to *Dicistroviridae*  
183 family, with 90 and 36 contigs from Quantuck Bay and Narragansett Bay, respectively (Fig. 3).  
184 Interestingly, the only dsRNA viruses detected in both locations were similar to viruses in the  
185 *Totiviridae*, *Partitiviridae* and *Hypoviridae* family, which are all known viruses of fungi (21).  
186 These viruses may be infecting fungi that are parasitic on algae, as have been proposed recently  
187 for samples collected in the Laurentian Great Lakes (25). While some ssDNA virus contigs from  
188 Quantuck Bay clustered with the *Nanoviridae* family, others from both locations did not form  
189 any definitive cluster with known circular DNA viruses, thus potentially representing previously  
190 undescribed circular DNA virus groups in the ocean (Fig. S2). No (-) ssRNA viral contigs were  
191 detected.

192 To assess how the activity of virus groups changed over time, we measured the proportion of  
193 reads that mapped to different virus groups for each library. In Narragansett Bay, the majority of  
194 the virus reads originated from the unclassified marine Picornavirales and the *Dicistroviridae*,  
195 *Secoviridae* and *Picornaviridae* families across all the time points (Fig. S3). The unclassified  
196 marine Picornavirales group recruited from ~68% (NB-S1) to ~98% (NB-S5) of the non-  
197 NCLDV viral reads (Fig. S3). In Quantuck Bay, reads from both unclassified marine  
198 Picornavirales and ssDNA viruses dominated libraries during the first two time points (Fig. S4).  
199 However, a shift in the proportional abundance of virus reads was observed on the third day  
200 (QB-S3), when the unclassified marine Picornavirales became dominant (93% of the non-  
201 NCLDV virus transcripts, Fig. S4). Overall, 2.4% of the entire QB-S3 library (~4.3 million  
202 fragments) mapped to these unclassified Picornavirales, relative to 0.043% and 0.027% of reads  
203 for QB-S1 and QB-S2, respectively. This indicated a striking increase, concordant with the  
204 decline of the brown tide bloom. In parallel nutrient amended samples derived from water  
205 collected on the same date as QB-S3, the unclassified Picornavirales transcripts also increased  
206 by an order of magnitude relative to QB-S1 and QB-S2 (Data not shown) validating the increase  
207 in viral reads in QB-S3. Phylogenetic analysis confirmed these ssRNA viral contigs were



208 consistent with the unclassified Picornavirales group (Fig. 3A). *Aureococcus* blooms are not  
209 mono-specific – they include diatoms, dinoflagellates and high densities ( $\sim 10^4$  cells/ml) of  
210 heterotrophic protists (26). The striking increase in the unclassified Picornavirales could be  
211 related to infection of a host that co-occurs with *Aureococcus* and the potential shift in  
212 competition that might occur during *Aureococcus* bloom decline. These observations suggest a  
213 broad ecological role for viral infection during phytoplankton bloom decline, which would not  
214 have been resolved with targeted studies of AaV or metagenomic approaches. Taken together,  
215 the apparent dynamics and abundance of this unclassified Picornavirales suggest this group is a  
216 major component of the marine viroplankton, and strengthens previous observations that the  
217 Picornavirales phylogenetically distinct from the established families can be dominant members  
218 in different marine environments (12, 27). Owing to their small size, detection and  
219 quantification of RNA viruses and ssDNA viruses pose significant technical challenges (27). Our  
220 results, however, clearly point to the power of screening metatranscriptomes for the simultaneous  
221 analysis of the dynamics of a large cross-section DNA and RNA viruses.

222 Eighteen of the assembled contigs (9 from each site) were  $>7,000$  bp and had best hits to  
223 different Picornavirales members. Phylogenetic analysis based on the RNA-dependent RNA  
224 polymerase (RdRp) gene and feature analysis of existing (+) ssRNA virus genomes suggested  
225 that these contigs are complete or near-complete Picornavirales genomes (Fig. 4). Sixteen of  
226 these genomes were dicistronic – they harbored two ORFs coding for structural and non-  
227 structural proteins, while the remaining 2 were monocistronic (Fig. 4), revealing differences in  
228 genome architecture among group members (Fig. 4). Remarkably, one virus (N\_001) had a  
229 reverse orientation of the genes with the first ORF encoding for the structural protein, which is  
230 unusual for dicistronic Picornavirales (28). In addition, a glycosyl transferase domain was found  
231 in N\_137 (Fig. 4). To our knowledge, the presence of glycosyltransferase domains has only been  
232 reported in members of the *Endornaviridae* family dsRNA viruses (29).

233 We also tracked the dynamics of these *de novo* assembled genomes by mapping the data  
234 collected over spatiotemporal gradients. All the (+) ssRNA virus genomes from Quantuck Bay  
235 samples showed higher relative expression during bloom decline (QB-S3) compared to the time  
236 points corresponding to the bloom peak (Fig. 4, panel C). N\_001, a candidate virus from  
237 Narragansett Bay, was not present in the first three time points. Its expression was only observed  
238 during the fourth sampling point, which was followed by a dramatic increase during the last

239 sampling point, when it recruited ~0.55% of the reads from the entire library (Fig. 4, panel C).  
240 The closest known phylogenetic relative of N\_001 is a virus infecting diatom *Asterionellopsis*  
241 *glacialis* (Fig. 4), suggesting the putative host may be a diatom. Narragansett Bay was  
242 experiencing a spring diatom bloom during the sampling period with “boom–bust” abundance  
243 cycles in the relative abundance of putative diatom hosts (30), consistent with these observed  
244 viral dynamics.

245

### 246 **Who is infecting whom? Resolving virus-host relationships using metatranscriptomics**

247 This study presented the opportunity to evaluate potential relationships among diverse single-  
248 celled eukaryotes and their pathogens, with the established *AaV-Aureococcus* association acting  
249 as a *de facto* internal standard. Transcripts from DNA viruses must originate within the host  
250 cells, and thereby for a particular host-virus pair, a significant and strong positive correlation is  
251 to be expected for gene expression. Building on this idea, host gene expression of at least a  
252 subset of the host’s genome is a prerequisite to observe gene expression of a virus specific to that  
253 host, as evidenced by transcriptomic landscape of host-virus dynamics in culture (17, 31) and  
254 induced blooms in mesocosms (32). To expand our data, we also took advantage of concurrent  
255 nutrient amendment studies in mesocosms (see Methods and Materials) which provided  
256 additional samples for our analyses.

257 We inspected statistical co-occurrences among the contigs containing virus and eukaryote-  
258 specific marker genes based on their expression patterns. Since the polyadenylated-selected  
259 metatranscriptomes are largely depleted of ribosomal RNA marker genes, we employed  
260 functional genes suitable for phylogenetic analysis. Expression of MCP (dsDNA NCLDVs),  
261 RNA dependent RNA polymerase (RdRP) (RNA viruses) and viral replicase (ssDNA viruses)  
262 were compared to the functional eukaryotic marker gene RNA polymerase II large subunit  
263 (RPB1, Fig. S5), a candidate gene to resolve the phylogenetic history of different eukaryotic  
264 lineages (33, 34). Hierarchical clustering of a Pearson’s correlation matrix followed by  
265 SIMPROF analysis (35) was used to detect statistically distinct clusters which contained both  
266 virus and eukaryotic marker genes that could be classified into phylogenetic groups. This  
267 analysis produced several statistically distinct clusters harboring both viral and eukaryotic  
268 contigs (Fig. 5). A single cluster (Fig. 5A(ii)) harbored both *AaV* and *Aureococcus*, validating

269 that established ecologically relevant relationships between viruses and their hosts can be  
270 retrieved using this approach.

271 Close inspection revealed other interesting relationships among the coexisting eukaryotic and  
272 viral components. Cluster A(ii), while containing both *Aureococcus* and AaV, also contained  
273 another *Mimiviridae* member, several ssDNA and (+) ssRNA viral contigs along with  
274 eukaryotes belonging to prasinophyceae and pelagophyceae (Fig. 5). The possibility of  
275 *Aureococcus* being infected by more than one virus type cannot be discounted (and indeed is  
276 perhaps likely). Moreover, the potential for AaV to infect closely related pelagophytes remains a  
277 possibility (although this has not been seen in lab studies) (36). One cluster, A(i), which  
278 contained both a *Phycodna*- and a *Mimiviridae* member, also included a RPB1 contig  
279 phylogenetically placed in the Cercozoa group (Fig.5). Although no cercozoan host-NCLDV  
280 pairs currently exist in culture, a recent study showed integration of NCLDV genes in the  
281 genome of a cercozoan *Bigelowella natans* (37). This integrated NCLDV in the *B. natans*  
282 genome potentially belongs to the *Phycodnaviridae*, as revealed by phylogenetic analysis of the  
283 MCP gene.

284 Similar clusters of phylogenetically distinct eukaryotes and viruses were also found in  
285 Narragansett Bay. For example, cluster B(iii) contained a *Mimiviridae* and several ssRNA viral  
286 contigs connected to choanomonada, stramenopile, diatom and dinoflagellate members (Fig. 5).  
287 The majority of the eukaryotic contigs belonged to diatoms and dinoflagellates in the  
288 Narragansett Bay samples, which reflects the composition of single-celled eukaryotes in this site  
289 (30). A large number of contigs having phylogenetic affinity to choanomonada were found in  
290 both Quantuck Bay and Narragansett Bay locations and were represented in several of the  
291 representative SIMPROF clusters (Fig. 5). While larger networks of viruses and eukaryotes also  
292 existed, clusters with fewer members revealed more specific relationships. For example, cluster  
293 B(xiv) contained one *Mimiviridae*, one jakobida (heterotrophic flagellate) and several diatom  
294 contigs (Fig. 5). To date the obligate heterotrophs known to be infected by *Mimiviridae* members  
295 are *Cafeteria roenbergensis* (38), *Acanthamoeba* (39) and *Vermaamoeba* spp (40). Additionally;  
296 Cluster B(xviii) harbored a ssDNA virus, a stramenopile and a choanomonada member, while  
297 cluster B(xxii) revealed a one-to-one relationship between a *Mimiviridae* and a dinoflagellate  
298 (Fig. 5). Only one dinoflagellate – *Heterocapsa circularisquama*, has been shown to be infected  
299 by a NCLDV (41), so this potential host-virus pair is of particular note. Cluster B(x) and B(xvii)

300 consisted of *Mimiviridae*, diatoms and ssRNA viruses. No diatom is yet known to be infected by  
301 a NCLDV, although a large number of ssRNA viral contigs in our study are phylogenetically  
302 close to diatom-infecting RNA viruses in the unclassified marine picornavirales group (Fig. 3). A  
303 number of clusters (*e.g.*, B(xii)) were enriched with both ssRNA virus and diatom contigs.  
304 These relationships between ssRNA viruses and the eukaryotes needs to be interpreted with  
305 caution, however, as these contigs might originate both from free virus particles and/or viruses  
306 within hosts.

307 Several clusters also contained fungal contigs along with other eukaryotes – pointing to the  
308 possibility of broad parasitic relationships with phytoplankton and other single-celled  
309 eukaryotes. The AaV-*Aureococcus* cluster A(ii) harbored a fungal contig and a *Barnaviridae*  
310 member – a virus family with fungi as the only known hosts (Fig. 5) (21). Several other clusters,  
311 *e.g.*, A(iii) and B(vii) also contained fungal contigs. While such observations are not definitive,  
312 they point to the existence of parasitic relationships resulting in complicated ecological  
313 interactions involving unicellular eukaryotes, fungi and fungal viruses in marine ecosystems  
314 (Edgar *et al.*, 2016).

315 Increased sample resolution in the future will resolve more statistically robust relationships,  
316 which can further narrow potential interacting partners. One limitation of reference independent  
317 assembly of high throughput data is fragmented contigs originating from same transcript – which  
318 is illustrated by two *Aureococcus* specific RPB1 contigs in cluster A(ii) that originated from a  
319 single coding sequence. Increased sequencing, along with the continued development of  
320 assembly tools will provide better resolution to these relationships. These limitations  
321 notwithstanding, the analysis provides a ‘proof-of-principle’ for inferring the complex  
322 relationships among diverse unicellular eukaryotes and their viruses using metatranscriptome  
323 data.

324

## 325 **Conclusion:**

326 In this study we demonstrate how metatranscriptomics can provide a unique view of the marine  
327 virosphere by simultaneously detecting multiple viral infections across the landscape of the  
328 eukaryotic plankton within an ecosystem setting. This effort can largely overcome previous  
329 technical limitations involved in the study of different viral groups, owing to their size range and

330 genome type, within the same sample. In the last two decades we have learned much regarding  
331 the diversity and dynamics of the phages in the ocean, but the eukaryotic virosphere has  
332 remained elusive, with little known about who is infecting whom in the environment (42). As  
333 demonstrated in our study, analyzing the vast wealth of information captured by  
334 metatranscriptomics, in a statistical framework, can be an important step towards answering this  
335 vital question.

336

337 **Acknowledgements:**

338 The authors thank Gary LeClerc and Eric Gann for assistance in data collection and analyses.  
339 This work was funded by grants from the Gordon and Betty Moore Foundation (Grant EMS4971  
340 to SWW), The National Science Foundation (OCE-1061352) and the *Kenneth & Blaire*  
341 *Mossman Endowment* to the University of Tennessee (SWW). Additional funding was provided  
342 by the NOAA ECOHAB Program (NA15NOS4780199 to STD and CJG).

343

344

345

346

347 **Materials and Methods:**

348

349 **Experimental design**

350 **Quantuck Bay:** Samples were collected from a brown tide bloom in Quantuck Bay (Latitude =  
351 40.806395; Longitude = -72.621002), NY that occurred from late May to early July, 2011,  
352 covering the initiation, peak and demise of the bloom. Samples collected on June 14 (BT-S1) and  
353 June 16 (BT-S2) represented the peak of the bloom, while sample collected on June 22 (BT-S3)  
354 represented the initial phase of bloom decline. *Aureococcus* cells were counted from  
355 glutaraldehyde (1% final v/v) fixed whole water samples using an enzyme-linked  
356 immunosorbent assay (ELISA) with a monoclonal antibody as described previously (43). *In situ*  
357 samples from June 22<sup>nd</sup> (3rd sampling point) was also used to carry out nutrient amendment  
358 experiments. Briefly, bottles were filled with natural sea water from the bloom and were  
359 amended with 25  $\mu$ M ammonium only (+N), 4  $\mu$ M phosphate only (+P), and 25  $\mu$ M ammonium  
360 and 4  $\mu$ M phosphate (+N&P) in triplicate. Three additional bottles with no nutrient addition were  
361 used as control. The samples were then incubated for 24 hours in a floating chamber at 0.5 m in  
362 eastern Shinnecock Bay at the Stony Brook - Southampton Marine Science Center under one  
363 layer of neutral density cover to mimic the light and temperature levels of Quantuck Bay.  
364 Samples for *Aureococcus* cell density measurement and total RNA extraction were collected at  
365 T=0 and T=24 hours. Approximately 25 ml of natural seawater from each of the *in situ* and  
366 nutrient amendment samples were pre-filtered through 5  $\mu$ m polycarbonate (PC) filters and  
367 collected on 0.2  $\mu$ m PC filters. The samples were flash frozen immediately after filtration and  
368 transferred to -80<sup>o</sup> C. Prior RNA extraction, CTAB buffer (Teknova, CA, USA) amended by  
369 polyvinylpyrrolidone (1% mass/vol) was added to each of the samples.

370

371 **Narragansett Bay:** The details sampling procedure is described in Alexander et al (30).  
372 Briefly, samples were collected from a long term sampling site in Narragansett Bay (41°34'12''  
373 N, 71°23'24'' W) during 2012 on May 16 (NB-S1), May 21 (NB-S2), May 30 (NB-S3), June 4  
374 (NB-S4) and June 8 (NB-S5). Sample collection and processing was completed within 0830 and  
375 0900 local time to reduce the influence of diel signals. 2.0 L of water from each sample was  
376 filtered on 5.0- $\mu$ m pore size PC filters using a peristaltic pump. The filters were snap frozen in

377 liquid nitrogen and stored at  $-80^{\circ}\text{C}$  until RNA extraction. Water collected along with NB-S3 was  
378 also used for nutrient amendment experiments. For this, triplicate 2.5L bottles were filled with  
379 water pre-filtered through a 200- $\mu\text{m}$  mesh and amended with specific nutrients to create +N, +P,  
380 -N, -P treatments alongside an ambient control. The +N and +P treatments were designed to  
381 eliminate nitrogen and phosphate stress signals, whereas the -N and -P treatments were designed  
382 to drive the treatments towards each limitation respectively by skewing the nutrient ratios  
383 (Alexander *et al.*, 2015). N and P amendment concentrations were ~10-fold the seasonal average  
384 N and P concentrations measured at the station II in the surface waters of Narragansett Bay. The  
385 +P and +N amendment contained 3  $\mu\text{M}$  phosphate and 10  $\mu\text{M}$  nitrate, respectively. The -P  
386 amendment contained 10  $\mu\text{M}$  nitrate, 68  $\mu\text{M}$  silica, 4.6  $\mu\text{M}$  iron and f/5 vitamins. The -N  
387 treatment was amended with 3  $\mu\text{M}$  phosphate, 68  $\mu\text{M}$  silica, 4.6  $\mu\text{M}$  iron and f/5 vitamins. The  
388 f/5 media ratios (44) were followed for silica, iron and vitamin amendments. Bottles were  
389 incubated for 48 hours in a flow-through incubator at ambient temperature and  
390 photosynthetically active radiation. After the end of the incubation, treated and control samples  
391 were filtered and stored for RNA extraction in the same manner for the *in situ* samples.

392

### 393 **RNA extraction and sequencing:**

394 Quantuck Bay samples were extracted using the UltraClean® Plant RNA Isolation Kit (MO BIO  
395 Laboratories, CA, USA) according to manufacturer's protocol. RNA samples were quantified  
396 spectrophotometrically and were sequenced in the Columbia Sequencing Center (NY, USA)  
397 using Illumina™ HiSeq™ platform with poly-A enrichment at a depth of ~50 million 100bp  
398 single end reads. Two more replicate samples were sequenced from June 22 (QB-S3) at a depth  
399 of 100 million reads. For the present study, these biological replicates from QB-S3 were pooled  
400 together prior to further analysis. The field sequence data reported in this paper have been  
401 deposited in the National Center for Biotechnology Information Sequence Read Archive,  
402 [www.ncbi.nlm.nih.gov/sra](http://www.ncbi.nlm.nih.gov/sra) (Narragansett Bay accession no. SRP055134; Quantuck Bay accession  
403 no. SRP072764).

404 For Narragansett Bay, replicate filters from each treatment and *in situ* samples were pooled,  
405 representing 6 L of water for each sample. RNA was extracted using RNeasy Mini Kit (Qiagen,  
406 Hilden, Germany) according to a modified yeast RNA extraction protocol. Briefly, lysis buffer



407 and RNA-clean zircon beads were added to the filter. Samples were then vortexed for 1 min,  
408 placed on ice for 30 s, and then vortexed again for 1 min. The resulting RNA was eluted in water  
409 and possible DNA contamination was removed using a TURBO DNA-free Kit (Thermo Fisher  
410 Scientific, MA, USA). RNA from each triplicate was pooled by sample or treatment. >1,000 ng  
411 RNA from each sample then went through a poly-A selection using oligo-dT beads followed by  
412 library preparation with TruSeq RNA Prep Kit (Illumina, CA, USA). The samples were  
413 sequenced with an Illumina™ HiSeq2000™ at the Columbia University Genome Center to  
414 produce ~60 million; 100bp paired-end reads per sample.

#### 415 **Read assembly and screening for virus and eukaryote specific contigs:**

416 Sequence reads from both locations were quality trimmed (stringent trimming (quality score  
417  $\leq 0.03$ ), No ‘N’s allowed, 70bp size cutoff) in CLC genomics workbench 8.0 (Qiagen, Hilden,  
418 Germany). The data from all the Quantuck Bay samples were combined and assembled together,  
419 and a similar assembly procedure was also followed for the Narragansett Bay specific samples.  
420 This resulted in 2,455,926 contigs for Quantuck Bay and 9,525,233 contigs for Narragansett Bay  
421 at a 100bp size cut-off.

422 For selecting contigs specific to Major Capsid proteins of NCLDV, a HMM profile was created  
423 after aligning the MCP sequences from complete giant virus genomes and several reported MCP  
424 genes available in NCBI. The HMM profile was queried against the translated contig libraries to  
425 select the putative MCP candidate contigs using HMMER (45). For selecting eukaryotic RPB1  
426 contigs, HMM profile specific to domain “RPB1-C-term (NCBI CDD ID: cd02584)” and  
427 “RPB1-N-term (NCBI CDD ID: cd02733)” was used to query the contig libraries. All the MCP  
428 and RPB1 candidate contigs detected in this manner were queried against NCBI Refseq database  
429 and only contigs with first BLASTx hits (e-value cut-off  $\leq 10^{-3}$ ) to MCP and RPB1 were kept for  
430 further analysis.

431 To detect contigs originating from viruses other than NCLDVs, we combined the proteins  
432 derived from all the viruses having algal, fungal and protozan hosts available on NCBI database.  
433 This protein database was queried against the contig libraries using tBLASTn with an e-value  
434 cut-off of  $\leq 10^{-3}$ . All the candidate contigs screened by this procedure were then queried against  
435 NCBI Refseq database using BLASTx. Only contigs having topmost hits to different viruses  
436 were kept for further analysis. All these contigs had best hits to diverse eukaryotic viruses- no

437 hits to prokaryotic viruses were detected. This is probably due to the fact that the samples were  
438 poly-A selected prior to sequencing. These virus contigs were binned into distinct viral groups  
439 according to their best BLASTx hits. Percentage of reads recruited to individual viral groups was  
440 calculated for determining proportional abundance of different viral groups over different time  
441 points. For detailed phylogenetic analysis of ssRNA and ssDNA viruses, subset of these contigs  
442 harboring RdRP (pfam id: PF05183) and viral replicase (pfam: PF03090) motif were selected  
443 using HMM profile specific to these motifs.

444

#### 445 **Genomic and phylogenetic analysis:**

446 Reference sequences for MCP (giant viruses), RdRP (RNA viruses), viral replicase (ssDNA  
447 virus) and RPB1 (eukaryotes) were downloaded from NCBI Refseq database. A number of  
448 RPB1 sequences representing several eukaryotic groups were also collected from Marine  
449 Microbial Eukaryotes Transcriptome Sequencing Project (MMETSP) (46) peptide collections,  
450 for which no representative genomes are available. The reference sequences were aligned in  
451 MEGA 6.0 (47). Maximum likelihood phylogenetic trees were constructed in PhyML (48) with  
452 LG model, gamma shape parameter and frequency type estimated from the data. aLRT SH-like  
453 statistic was calculated for branch support. The eukaryotic classification scheme by Adl *et al.*  
454 (49) was followed. Selected contigs were translated to amino acid sequences and were placed on  
455 the reference trees in a maximum likelihood framework using pplacer (50). The placement files  
456 were converted to trees with pendant edges showing the best placement of the contigs using  
457 ‘guppy’ tool of pplacer. The placement trees were visualized and annotated using iTOL interface  
458 (51).

459 ORFs were predicted on the complete or near-complete Picornavirales genomes using CLC  
460 genomic workbench 8.0 ([www.clcbio.com](http://www.clcbio.com)). The genome annotation with predicted features was  
461 assisted by pfam (52) and Conserved Domain Database (CDD) (53) search. The genomes are  
462 submitted in the NCBI database under the accession numbers: KY286099 - KY286107 and  
463 KY130489 - KY130497.

464

#### 465 **Statistical analysis:**

466 Quality trimmed reads were mapped to the selected viral and eukaryotic contigs from  
467 individual read libraries with high stringency (97% identify, 70% length fraction matching) in  
468 CLC genomics workbench 8.0. The read mapping values were normalized by library size and  
469 length and expressed as rarefied counts per kilobase (RCK). RCK values of viral and eukaryotic  
470 contigs >225 base pairs were converted into matrices separately for Quantuck Bay and  
471 Narragansett Bay datasets, which included mapping statistics from both *in situ* and nutrient  
472 amendment libraries. Group averaged hierarchical clustering was performed on these matrices  
473 using Pearson's correlation coefficient in PRIMER 7.0. SIMPROF test (35) was applied on the  
474 clusters with 5% significance level and 1000 permutations to identify statistically distinct  
475 clusters. Selected clusters were visualized and annotated in Cytoscape 3.0 (54).

476

477

478 **References:**

- 479 1. Breitbart M (2012) Marine viruses: truth or dare. *Ann Rev Mar Sci* 4:425-448.  
480 2. Brum JR, *et al.* (2015) Patterns and ecological drivers of ocean viral communities.  
481 *Science* 348(6237).  
482 3. Weitz JS & Wilhelm SW (2012) Ocean viruses and their effects on microbial  
483 communities and biogeochemical cycles. *F1000 Biol Rep* 4:17.  
484 4. Short SM (2012) The ecology of viruses that infect eukaryotic algae. *Environ Microbiol*  
485 14(9):2253-2271.  
486 5. Ogata H, Monier A, & Claverie J-M (2011) Distribution of Giant Viruses in Marine  
487 Environments. *Global Change: Mankind-Marine Environment Interactions*, eds Ceccaldi  
488 H-J, Dekeyser I, Girault M, & Stora G (Springer Netherlands), pp 157-162.  
489 6. Lehahn Y, *et al.* (2014) Decoupling physical from biological processes to assess the  
490 impact of viruses on a mesoscale algal bloom. *Curr Biol* 24(17): 2041-2046  
491 7. Gastrich M, *et al.* (2004) Viruses as potential regulators of regional brown tide blooms  
492 caused by the alga, *Aureococcus anophagefferens*. *Estuaries* 27(1):112-119.  
493 8. Abergel C, Legendre M, & Claverie J-M (2015) The rapidly expanding universe of giant  
494 viruses: Mimivirus, Pandoravirus, Pithovirus and Mollivirus. *FEMS Microbiol Rev*  
495 39(6):779-796  
496 9. Lang AS, Rise ML, Culley AI, & Steward GF (2009) RNA viruses in the sea. *FEMS*  
497 *Microbiol Rev* 39 33(2):295-323.  
498 10. Tai V, *et al.* (2003) Characterization of HaRNAV, a single-stranded RNA virus causing  
499 lysis of *Heterosigma akashiwo* (Raphidophyceae). *J Phycol* 39(2):343-352.  
500 11. Labonte JM & Suttle CA (2013) Previously unknown and highly divergent ssDNA  
501 viruses populate the oceans. *ISME J* 7(11):2169-2177.  
502 12. Steward GF, *et al.* (2013) Are we missing half of the viruses in the ocean? *ISME J*  
503 7(3):672-679.  
504 13. Gobler CJ, *et al.* (2011) Niche of harmful alga *Aureococcus anophagefferens* revealed  
505 through ecogenomics. *Proc Natl Acad Sci USA*. 108 (11):4352-4357  
506 14. Gastrich MD, Anderson OR, Benmayor SS, & Cospér EM (1998) Ultrastructural analysis  
507 of viral infection in the brown-tide alga, *Aureococcus anophagefferens* (Pelagophyceae).  
508 *Phycologia* 37(4):300-306.  
509 15. Gastrich MD, Leigh-Bell JA, Gobler C, Anderson OR, & Wilhelm SW (2004) Viruses as  
510 potential regulators of regional brown tide blooms caused by the alga, *Aureococcus*  
511 *anophagefferens*: a comparison of bloom years 1999-2000 and 2002. *Estuaries*  
512 27(1):112-119.  
513 16. Yutin N, Wolf Y, Raoult D, & Koonin E (2009) Eukaryotic large nucleo-cytoplasmic  
514 DNA viruses: Clusters of orthologous genes and reconstruction of viral genome  
515 evolution. *Virol J* 6(1):223.  
516 17. Legendre M, *et al.* (2010) mRNA deep sequencing reveals 75 new genes and a complex  
517 transcriptional landscape in Mimivirus. *Genome Res* 20(5):664-674.  
518 18. Moniruzzaman M, *et al.* (2016) Diversity and dynamics of algal Megaviridae members  
519 during a harmful brown tide caused by the pelagophyte, *Aureococcus anophagefferens*.  
520 *FEMS Microbiol Ecol*. 92: fiw058  
521 19. Brussaard CPD (2004) Viral control of phytoplankton populations—a review. *J Euk*  
522 *Microbiol* 51(2):125-138.

- 523 20. Fløge SA (2014) Virus infections of the eukaryotic marine microbes. PhD Dissertation  
524 (The University of Maine, Maine).
- 525 21. Wilson WH, Van Etten, J.L., Schroeder, D.C., Nagasaki, K., Brussaard, C., Bratbak, G.,  
526 Suttle, C. (2012) Virus Taxonomy: Ninth Report of the International Committee on  
527 Taxonomy of Viruses. ed Andrew M.Q. King EL, Michael J. Adams , Eric B. Carsten  
528 (Elsevier London, UK), pp 10-20.
- 529 22. Shatkin AJ (1974) Animal RNA viruses: genome structure and function. *Ann Rev*  
530 *Biochem* 43(1):643-665.
- 531 23. Culley AI, Lang AS, & Suttle CA (2003) High diversity of unknown picorna-like viruses  
532 in the sea. *Nature* 424(6952):1054-1057.
- 533 24. Takao Y, Mise K, Nagasaki K, Okuno T, & Honda D (2006) Complete nucleotide  
534 sequence and genome organization of a single-stranded RNA virus infecting the marine  
535 fungoid protist *Schizochytrium* sp. *J Gen Virol* 87(Pt 3):723-733.
- 536 25. Edgar RE, *et al.* (2016) Adaptations to photoautotrophy associated with seasonal ice  
537 cover in a large lake revealed by metatranscriptome analysis of a diatom bloom. *J Gt*  
538 *Lakes Res* 42 (5):1007-1015.
- 539 26. Sieracki ME, *et al.* (2004) Pico- and nanoplankton dynamics during bloom initiation of  
540 *Aureococcus* in a Long Island, NY bay. *Harmful Algae* 3(4):459-470.
- 541 27. Miranda JA, Culley AI, Schvarcz CR, & Steward GF (2016) RNA viruses as major  
542 contributors to Antarctic viroplankton. *Environ Microbiol.* 18(11):3714-3727.
- 543 28. Greninger AL & DeRisi JL (2015) Draft Genome Sequence of Laverivirus UC1, a  
544 Dicistrovirus-like RNA virus featuring an unusual genome organization. *Genome*  
545 *Announc* 3(4):e00656-00615.
- 546 29. Song D, Cho WK, Park S-H, Jo Y, & Kim K-H (2013) Evolution of and horizontal gene  
547 transfer in the *Endornavirus* genus. *PLoS ONE* 8(5):e64270.
- 548 30. Alexander H, Jenkins BD, Ryneerson TA, & Dyhrman ST (2015) Metatranscriptome  
549 analyses indicate resource partitioning between diatoms in the field. *Proc Natl Acad Sci*  
550 *USA* 112(17):E2182-E2190.
- 551 31. Rowe JM, *et al.* (2014) Global Analysis of *Chlorella variabilis* NC64A mRNA profiles  
552 during the early phase of *Paramecium bursaria* Chlorella Virus-1 infection. *PLoS ONE*  
553 9(3):e90988.
- 554 32. Pagarete A, *et al.* (2011) Unveiling the transcriptional features associated with  
555 coccolithovirus infection of natural *Emiliania huxleyi* blooms. *FEMS Microbiol Ecol*  
556 78(3):555-564.
- 557 33. Stiller JW & Hall BD (1997) The origin of red algae: implications for plastid evolution.  
558 *Proc Natl Acad Sci USA* 94(9):4520-4525.
- 559 34. Malik SB, *et al.* (2011) Phylogeny of parasitic parabasalia and free-living relatives  
560 inferred from conventional markers vs. Rpb1, a single-copy gene. *PLoS One*  
561 6(6):e20774.
- 562 35. Clarke KR, Somerfield PJ, & Gorley RN (2008) Testing of null hypotheses in  
563 exploratory community analyses: similarity profiles and biota-environment linkage. *J Exp*  
564 *Mar Biol Ecol* 366(1–2):56-69.
- 565 36. Gobler CJ, Anderson OR, Gastrich MD, & Wilhelm SW (2007) Ecological aspects of  
566 viral infection and lysis in the harmful brown tide alga *Aureococcus anophagefferens*.  
567 *Aquat Microb Ecol* 47(1):25-36.

- 568 37. Blanc G, Gallot-Lavallee L, & Maumus F (2015) Provirophages in the *Bigelowiella*  
569 genome bear testimony to past encounters with giant viruses. *Proc Natl Acad Sci U S A*  
570 112(38):E5318-5326.
- 571 38. Fischer MG, Allen MJ, Wilson WH, & Suttle CA (2010) Giant virus with a remarkable  
572 complement of genes infects marine zooplankton. *Proc Natl Acad Sci U S A*  
573 107(45):19508-19513.
- 574 39. Abrahão JS, *et al.* (2014) *Acanthamoeba polyphaga* mimivirus and other giant viruses: an  
575 open field to outstanding discoveries. *Virol J* 11:120.
- 576 40. Reteno DG, *et al.* (2015) Faustovirus, an asfarvirus-related new lineage of giant viruses  
577 infecting amoebae. *J Virol*.
- 578 41. Kenji T, Keizo N, Shigeru I, & Mineo Y (2001) Isolation of a virus infecting the novel  
579 shellfish-killing dinoflagellate *Heterocapsa circularisquama*. *Aquat Microb Ecol*  
580 23(2):103-111.
- 581 42. Caron DA, *et al.* (2017) Probing the evolution, ecology and physiology of marine protists  
582 using transcriptomics. *Nat Rev Micro* 15(1):6-20.
- 583 43. Koch F, Sañudo-Wilhelmy SA, Fisher NS, & Gobler CJ (2013) Effect of vitamins B1 and  
584 B12 on bloom dynamics of the harmful brown tide alga, *Aureococcus anophagefferens*  
585 (Pelagophyceae). *Limnol Oceanogr* 58(5):1761-1774.
- 586 44. Guillard RRL (1975) Culture of phytoplankton for feeding marine invertebrates. *Culture*  
587 *of Marine Invertebrate Animals*, ed Smith WL, Chanley, M.H. (Plenum Press, New York,  
588 USA), pp 26-60.
- 589 45. Eddy SR (2011) Accelerated profile HMM searches. *PLoS computational biology*  
590 7(10):e1002195.
- 591 46. Keeling PJ, *et al.* (2014) The marine microbial eukaryote transcriptome sequencing  
592 project (MMETSP): illuminating the functional diversity of eukaryotic life in the oceans  
593 through transcriptome sequencing. *PLoS Biol* 12(6):e1001889.
- 594 47. Tamura K, Stecher G, Peterson D, Filipinski A, & Kumar S (2013) MEGA6: molecular  
595 evolutionary genetics analysis version 6.0. *Mol Biol Evol* 30(12):2725-2729.
- 596 48. Guindon S, *et al.* (2010) New algorithms and methods to estimate maximum-likelihood  
597 phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* 59(3):307-321.
- 598 49. Adl SM, *et al.* (2005) The new higher level classification of eukaryotes with emphasis on  
599 the taxonomy of protists. *J Euk Microbiol* 52(5):399-451.
- 600 50. Matsen FA, Kodner RB, & Armbrust VE (2010) pplacer: linear time maximum-  
601 likelihood and Bayesian phylogenetic placement of sequences onto a fixed reference tree.  
602 *BMC Bioinformatics* 11(1):1-16.
- 603 51. Letunic I & Bork P (2016) Interactive tree of life (iTOL) v3: an online tool for the  
604 display and annotation of phylogenetic and other trees. *Nucl Acid Res* 44(W1):W242-  
605 W245.
- 606 52. Punta M, *et al.* (2012) The Pfam protein families database. *Nucl Acid Res* 40(D1):D290-  
607 D301.
- 608 53. Marchler-Bauer A, *et al.* (2013) CDD: conserved domains and protein three-dimensional  
609 structure. *Nucl Acid Res* 41(Database issue):D348-352.
- 610 54. Cline MS, *et al.* (2007) Integration of biological networks and gene expression data using  
611 Cytoscape. *Nature protocols* 2(10):2366-2382.

612 **Figures and tables:**

613 **Figure 1:** Abundance of 9 NCLDV core genes including major capsid protein (MCP) in terms of  
614 normalized read counts and number of contigs recovered (up to 100bp length) from Quantuck  
615 Bay (right panel) and Narragansett Bay (left panel). The box and whisker plots represent the  
616 range of the contig lengths with number of contigs recovered for each gene in brackets. The  
617 filled circles represent the rarefied abundances of each contig in each sample. No contigs could  
618 be detected from myristolyated envelope protein, a core NCLDV gene. Major capsid protein  
619 abundances are shown in red, while other core genes are presented in dark gray. Core genes are  
620 indicated on the X-axes as follows: A) A32 virion packaging ATPase, B) VLFT3 like  
621 transcription factor, C) Superfamily II helicase II, D) mRNA capping enzyme, E) D5  
622 helicase/primase, F) Ribonucleotide reductase small subunit, G) RNA polymerase large subunit,  
623 H) RNA polymerase small subunit, I) B-family DNA polymerase, J) Major capsid protein.

624

625 **Figure 2:** Phylogenetic placement of major capsid protein contigs from A) Quantuck Bay and B)  
626 Narragansett Bay on a reference tree of NCLDVs with icosahedral capsids. Node support (aLRT-  
627 SH statistic) >50% are shown as dark circles. Contigs upto 200bp are shown, with their  
628 expression level (rarefied read counts per kilobase – RCK) in individual samples as a heatmap on  
629 the outer rings. Notice that the placement trees contain both AaV reference sequence and the  
630 contig originating from AaV (marked with a black arrow). The reference sequences are shown in  
631 bold italic typeface. Abbreviations: MsV-Marseillevirus, LauV: Lausannevirus, Ws Irido:  
632 Weisenia iridescent virus, SG Irido: Singapore Grouper iridescent virus, He Asco: Heliothis  
633 virescens ascovirus, AsfV: Asfarvirus, EhV86: Emiliana huxleyi virus 86, HaV01: Heterosigma  
634 akashiwo virus 01, PBCV1: Paramacium bursaria Chlorella virus 1, ATCV 1: Acanthocystis  
635 turfacea Chlorella virus 1, BpV1: Bathycoccus parsinos virus 1, MpV12T: Micromonas pusilla  
636 virus 12T, OIV1: Ostreococcus lucimarinus virus 1, AaV: Aureococcus anophagefferens virus,  
637 CeV: Chrysochromulina ericina virus, PpV: Phaeocystis pouchetii virus, PgV: Phaeocystis  
638 globosa virus, PoV: Pyramimonas orientalis virus, Mega: Megavirus chilensis, Moumou:  
639 Moumouvirus goulette, Mimi: Mimivirus, CroV: Cafeteria roenbergensis virus.

640

641 **Figure 3:** Phylogenetic placement of (+)ssRNA virus contigs harboring RNA dependent RNA  
642 polymerase (RdRP) motifs from A) Quantuck Bay and B) Narragansett Bay on reference trees.  
643 Node support (aLRT-SH statistic) >50% are shown as dark circles. Contigs up to 225bp are  
644 shown, with their expression level (rarefied read counts per kilobase – RCK) in individual  
645 samples as a heatmap on the outer rings. The reference sequences are shown in bold italic  
646 typeface. Complete name and other details of the reference sequences are presented in  
647 Supplementary Dataset 2.

648

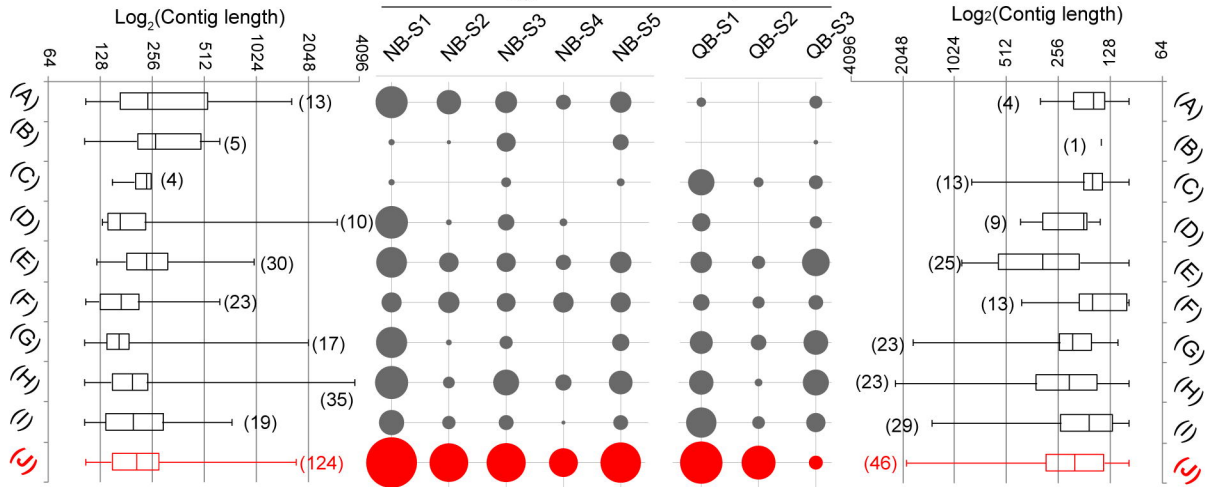
649 **Figure 4:** Complete or near-complete Picornavirales genomes recovered from both Quantuck  
650 Bay and Narragansett Bay study sites. Panel (A) shows the phylogenetic classification of these  
651 contigs in a topology-only maximum likelihood tree, with contigs from Quantuck Bay having  
652 prefix ‘Q\_’ and contigs from Narragansett Bay having prefix ‘N\_’. Panel (B) shows the genome  
653 architecture of these contigs with protein domains and putative CDSs. Panel (C) shows the  
654 expression level of these (rarefied read counts per kilobase – RCK) viruses in across different *in*  
655 *situ* samples.

656

657 **Figure 5:** Representative SIMPROF clusters containing both viral and eukaryotic members from  
658 A) Quantuck Bay and B) Narragansett Bay. Contigs are shown as nodes and the correlations as  
659 the connecting edges. Phylogenetic classifications of the contigs are shown in the bottom panel.  
660 *Aureococcus anophagefferens* (dark brown circles) and *Aureococcus anophagefferens* virus  
661 (AaV) (bright yellow square) are in cluster A(ii).



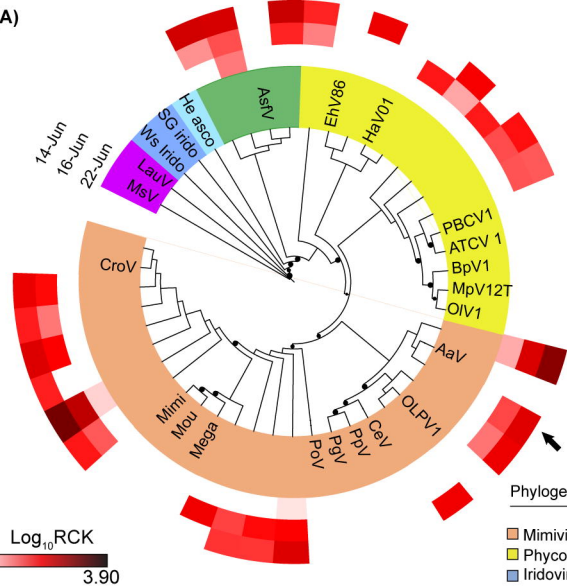
$\text{Log}_{10}(\text{Rarefied Read Counts})$



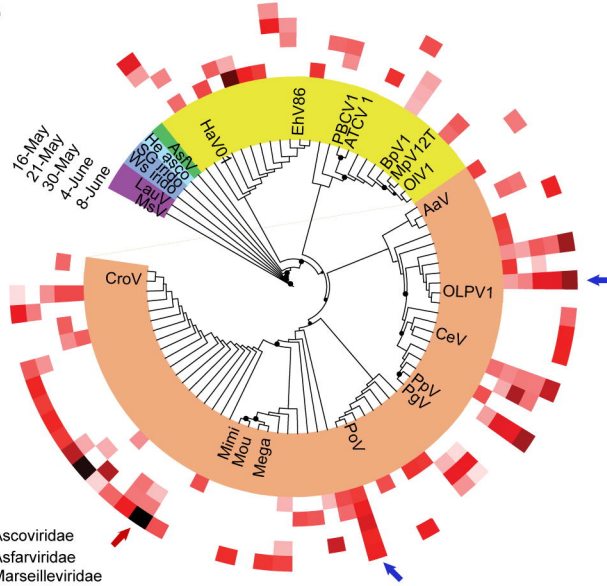
(i)

(ii)

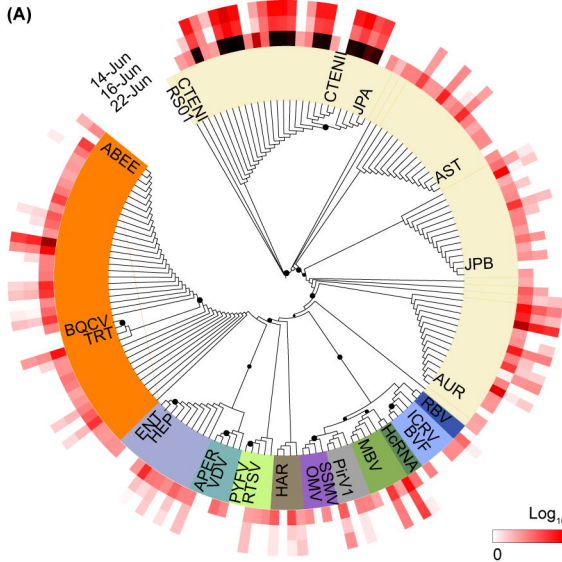
(A)



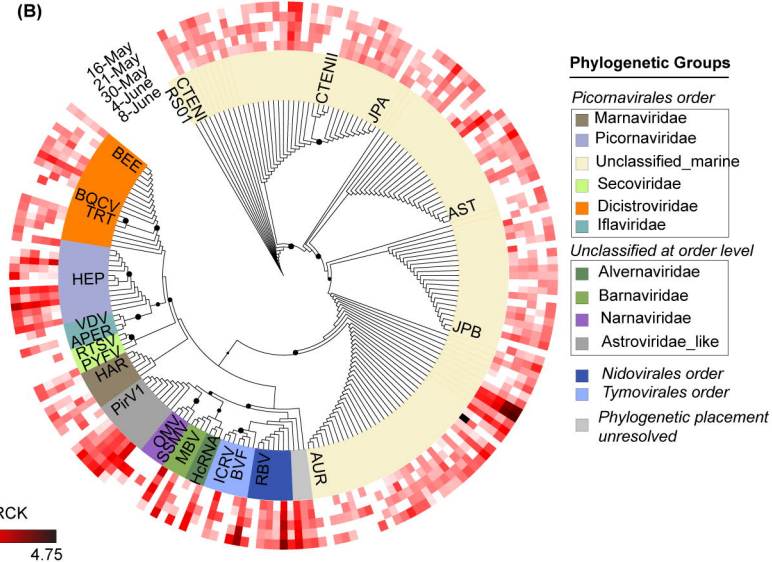
(B)

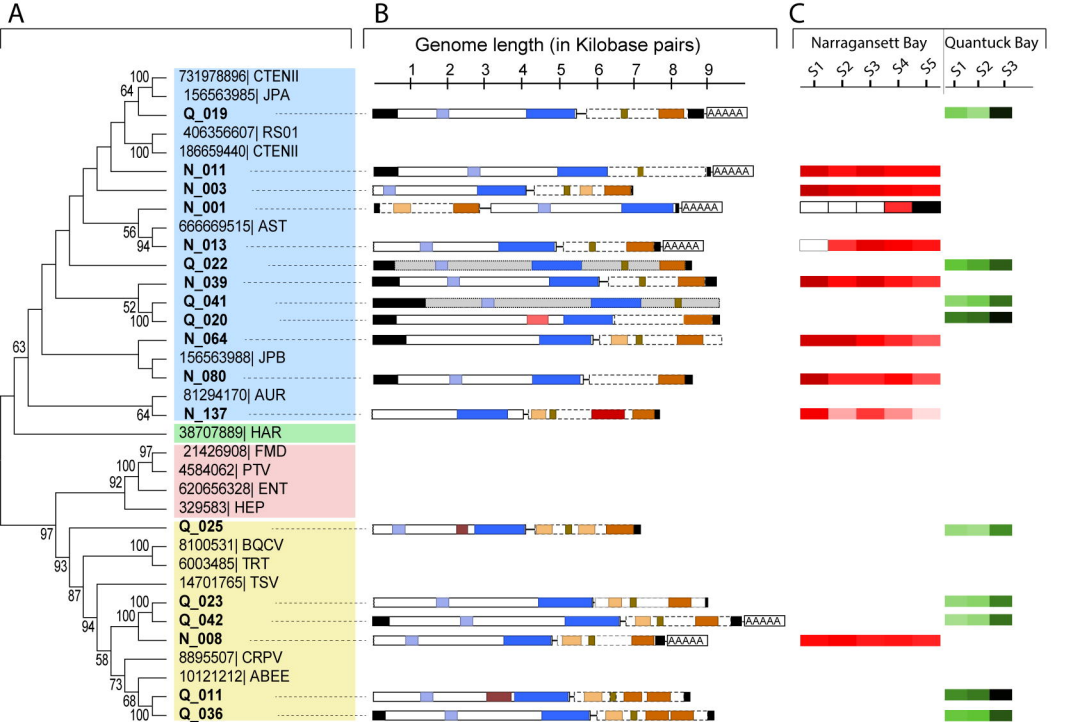


(A)

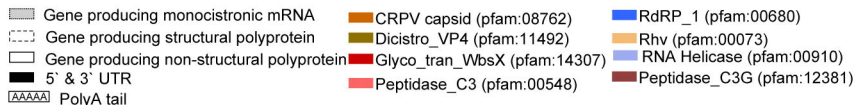


(B)



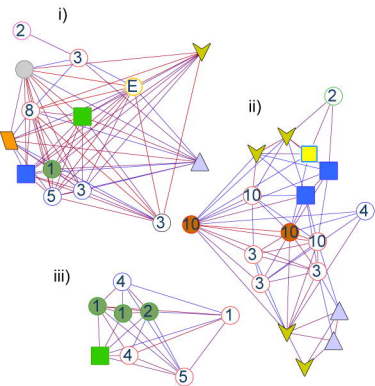
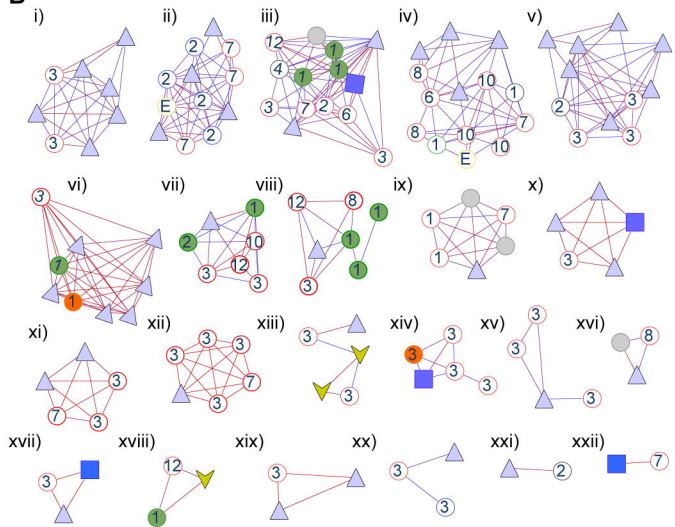


Annotations



LOG<sub>10</sub>(RCK)



**A****B**

■ Phycodnaviridae 
 ■ Mimiviridae 
 ■ Aureococcus virus (AaV) 
 ▲ (+)ssRNA virus 
 ▼ ssDNA virus 
 ■ dsRNA virus

Pearson correlation coefficient

