1	Interaction Dynamics and Virus-Host Range for Estuarine Actinophages captured by
2	epicPCR
3	Eric G. Sakowski <sup>1*</sup> , Keith Arora-Williams <sup>1</sup> , Funing Tian <sup>2</sup> , Ahmed A Zayed <sup>2</sup> , Olivier Zablocki <sup>2</sup> ,
4	Matthew B. Sullivan <sup>2,3</sup> , Sarah P. Preheim <sup>1*</sup>
5	
6	<sup>1</sup> Department of Environmental Health and Engineering, Johns Hopkins University, MD, USA
7	2 Department of Microbiology, The Ohio State University, OH, USA
8	3 Department of Civil, Environmental and Geodetic Engineering, The Ohio State University, OH
9	USA
10	
11	* Corresponding authors
12	
13	Viruses impact microbial diversity, gene flow, and function through virus-host
14	interactions. Though metagenomics surveys are rapidly cataloging viral diversity, methods are
15	needed to capture specific virus-host interactions in situ. We leveraged metagenomics and
16	repurposed emulsion paired isolation-concatenation PCR (epicPCR) to investigate viral diversity
17	and virus-host interactions in situ over time in an estuarine environment. The method fuses a
18	phage marker, the ribonucleotide reductase (RNR) gene, with the host 16S rRNA gene of
19	infected bacterial cells within emulsion droplets providing single-cell resolution for dozens of
20	samples. EpicPCR captured in situ virus-host interactions for viral clades with no closely related
21	database representatives. Abundant freshwater Actinobacteria lineages, in particular Rhodoluna,
22	were the most common hosts for these poorly characterized viruses, with interactions correlated
23	with environmental factors. Multiple methods used to identify virus-host interactions, including

epicPCR, identified different and largely non-overlapping interactions within the vast virus-host
interaction space. Tracking virus-host interaction dynamics also revealed that multi-host viruses
had significantly longer periods with observed virus-host interactions, while single-host viruses
were observed interacting with hosts at lower minimum abundances, suggesting more efficient
interactions. Capturing *in situ* interactions with epicPCR revealed environmental and ecological
factors shaping virus-host interactions, highlighting epicPCR as a valuable technique in viral
ecology.

31

Viruses impact the diversity and function of microbial communities from the open ocean<sup>1</sup> 32 and soils<sup>2</sup> to the human gut<sup>3</sup> and are major contributors to daily bacterial mortality<sup>4</sup>, carbon 33 export<sup>5</sup>, and biogeochemical cycling<sup>6</sup>. In the Chesapeake Bay, the largest estuary in the United 34 States, up to 20% of the bacterial community can be lysed by viruses per hour,<sup>7</sup> although 35 36 environmental factors (e.g. tidal mixing), could modulate viral production and viral-mediated mortality, as observed in other ecosystems<sup>8</sup>. However, cultivation<sup>9,10</sup> and theoretical models<sup>11,12</sup> 37 38 suggest viral pressure is unequally distributed across the microbial community, implying a subset 39 of virus-host infections contribute disproportionately to microbial mortality, influencing 40 microbial community diversity and biogeochemical cycling. Since individual virus-host 41 interactions (i.e. any physical association including attachment, injection of genetic material, 42 latent period, lysogeny) could be influenced by different ecological and environmental factors, 43 establishing connections between virus-host interactions and environmental factors will be the 44 first step in fine-tuning ecosystem models to better predict variability in viral productivity and 45 bacterial mortality.

46 Viral metagenomics has greatly expanded our knowledge of environmental viral 47 diversity; yet, linking viruses to their hosts remains a bottleneck for investigating the relative 48 ecosystem impact of individual virus-host pairs. Bioinformatics approaches applied to bacterial and viral shotgun metagenomes show promise for identifying hosts for viral communities<sup>13</sup>. 49 50 Techniques such as CRISPR spacer matches, sequence homology, and k-mer frequencies have been relatively successful, linking 35% of soil viral populations to putative hosts<sup>2</sup>. Floods of 51 oceanic virus data<sup>14</sup> and metagenome-assembled genomes (MAGs)<sup>15</sup> will likely improve *in silico* 52 53 host predictions for marine viruses. Still, only an estimated 10% of uncultivated bacterial lineages contain CRISPR-Cas systems<sup>16</sup>, limiting efforts to link viruses and hosts by CRISPR 54 55 spacer homology. Comparisons of mycobacteriophage suggest some viruses switch hosts too quickly for their genomes to evolve similar DNA signatures to their host<sup>17</sup>, thus constraining *in* 56 *silico* predictions via k-mer frequencies. Other techniques, like spot, plaque and liquid assays<sup>18</sup> 57 and viral tagging<sup>19,20</sup> require cultivation of the host, preventing investigations of viral 58 59 interactions for the majority of bacterial populations that cannot be readily cultivated. Cultivation-independent approaches like phageFISH<sup>21</sup>, polonies, or microfluidic digital 60 polymerase chain reaction<sup>22</sup> are important techniques in viral ecology, but require laborious 61 probe design and optimization, specialized expertise and equipment, or yield few positive 62 interactions. Other culture-independent techniques like proximity ligation<sup>23,24</sup> and single-cell 63 genomics<sup>25</sup> are exciting non-specific approaches for detecting virus-host interactions, but non-64 65 targeted approaches make consistently observing the same virus-host interactions over time 66 inefficient. A targeted and cost-effective approach capturing virus-host interactions in situ using 67 standard molecular biological equipment could fill a gap in current techniques that may reveal 68 ecological and environmental factors influencing specific virus-host interactions.

69	We sought to develop such a technique by repurposing emulsion paired isolation-
70	concatenation PCR $(epicPCR)^{26}$ to evaluate <i>in situ</i> virus-host interaction dynamics of
71	ecologically important bacterial populations. This analysis revealed abundant freshwater
72	Actinobacteria populations as the primary host of poorly characterized viruses sharing genetic
73	similarity with cyanosipho- and podoviruses. Additionally, environmental factors partially
74	explained the observed interactions between one abundant Actinobacteria population
75	(Rhodoluna) and viral clades. We also identified ecological differences in "interaction life-span"
76	and minimum host abundance between specialist and generalist viruses. These results
77	demonstrate that this approach complements existing techniques in viral ecology.
78	
79	Results
80	Chesapeake Bay viral populations are endemic, and seasonally dynamic
81	We used shotgun metagenomics to investigate temporal changes in viral populations and
82	virus-host interactions across two years in surface water samples in the Rhode River, a tidal
83	estuary of the Chesapeake Bay (Extended Data Fig. 1). In total, 9,392 viral populations
84	(approximately species level genotypes) were assembled from bacterial and viral shotgun
85	metagenomic libraries. Viral diversity in the Rhode River inlet (hereafter referred to as the
86	Chesapeake Bay) was compared to reference database representatives (RefSeq), Global Ocean
87	viromes [the Global Ocean Virome 2.0 (GOV 2.0) dataset <sup>14</sup> ] and U.S. East coast estuaries
88	(Damariscotta Salt Bay and Tampa Bay). Overall, Chesapeake Bay viral populations were
89	distinct from reference database representatives. Chesapeake Bay viral populations formed 473
90	clusters in a network analysis of shared genes <sup>27</sup> , but only 40 (8%) contained RefSeq
91	representatives. Over 2,000 Chesapeake Bay viral populations were not assigned a single cluster.

Additionally, only 5% and 1% of viral populations could be detected in other U.S. East coast
estuaries or GOV 2.0 datasets, respectively (Fig. 1). This highlights the potentially endemic
nature of Chesapeake Bay viral populations and is consistent with prior inferences of endemism
from cyanophage gene markers<sup>28</sup>. Endemism is likely attributable to the Bay's long residence
time (6-7 months) promoting distinctly estuarine microbial populations<sup>29</sup>.

Consistent with previous reports<sup>30,31</sup>, bacterial and viral communities displayed distinct 97 98 seasonal trends (Fig. 2). Winter viral communities from 2018 shared nearly three times more viral populations with a winter 2012 sample<sup>32</sup> than with spring and summer samples from the 99 100 same year (Fig. 2b, Extended Data Fig. 2a). Likewise, spring 2017 and 2018 samples shared half 101 of their viral populations with each other but fewer than 30% with other seasons in 2018 (Fig. 102 2b, Extended Data Fig. 2a), illustrating that Chesapeake Bay viral communities had significant 103 seasonal variability (p = 0.0003; Extended Data Fig. 2b) but no significant inter-annual 104 variability (Extended Data Fig. 2c).

105

#### 106 Actinobacteria are commonly identified hosts for observed viral populations

107 We used bioinformatics and experimental approaches to predict hosts for Chesapeake 108 Bay viral populations. We inferred hosts through similarity to reference sequences with the viral 109 marker gene ribonucleotide reductase (RNR), which is broadly distributed across aquatic dsDNA lytic viruses.<sup>33,34</sup> In total, 634 Chesapeake Bay viral populations contained RNR genes. RNR-110 111 containing populations accounted for  $5.9 \pm 0.8\%$  of viral populations and  $8.9 \pm 3.2\%$  of the viral 112 community in any given time point (Extended Data Fig. 3a,b) and captured the seasonal dynamics of the viral community ( $r_s = 0.93$ , p = 0; Extended Data Fig. 4), marking RNRs as a 113 114 good proxy of aquatic viral diversity in our study. Spring RNR genes were similar to

115 Cyanomyoviruses, Gammaproteobacteria siphoviruses, and Alphaproteobacteria podoviruses 116 (Extended Data Fig. 5a). The relative abundance of these myovirus and siphovirus populations 117 decreased in winter samples, while the relative abundance of podoviruses infecting 118 Alphaproteobacteria increased (Extended Data Fig. 5a). The summer sample was largely 119 composed of populations with no viral reference database representatives, and overall  $54 \pm 19\%$ 120 of the RNR sequences were not classified using this homology-based approach. However, even 121 'classifiable' RNR sequences were only distantly related to viral reference sequences (Extended 122 Data Fig. 5b), and homology to viral references does not necessarily indicate a common host. For 123 example, Pelagibacter phage HTVC008M and T4-like cyanophages share 71% RNR nucleotide 124 identity and would be classified together by our cutoffs; yet, they infect hosts from different 125 phyla. Therefore, gene homology may lack the resolution to predict hosts for environmental 126 viruses that are only distantly related to reference database representatives. 127 Cyanopodoviruses have been identified as dominant cyanophage populations in the Chesapeake Bay<sup>30</sup>; yet, cyanosipho-and-podoviral RNR genes from the Chesapeake Bay viral 128 129 metagenomes were only distantly related to reference cyanophage sequences (Extended Data 130 Fig. 5b). To investigate hosts associated with these RNR gene sequences, we modified an existing single-cell gene fusion protocol (epicPCR)<sup>26</sup> and applied it to our samples (Table S1). 131 132 EpicPCR takes advantage of the close physical proximity of phage and host DNA within infected 133 bacterial cells to fuse viral marker and host 16S ribosomal RNA (rRNA) genes within emulsion 134 droplets, providing single-cell level resolution (Fig. 3) even without a priori knowledge of 135 putative hosts. Auxiliary metabolic genes are an unusual choice for a viral marker given the possibility of horizontal gene transfer<sup>34</sup>, but reference cyanosipho-and-podoviral RNR genes 136 form a monophyletic clade (the Cyano SP clade<sup>35</sup>) distinct from other known viral and microbial 137

RNR genes, making Cyano SP RNR an ideal marker gene choice for epicPCR (Extended Data
Fig. 3c). Our primers specifically amplified RNR genes related to those of cyanosiphoviruses
and cyanopodoviruses (Table S2, Extended Data Fig. 3c). Furthermore, these amplicons were
dissimilar to microbial RNR genes from our bacterial shotgun metagenomes and UniRef<sup>36</sup>
representative sequences (Extended Data Fig. 3c), suggesting the primers specifically amplify
viral RNR genes.

144 After confirming primer specificity, we validated epicPCR specificity using uninfected 145 mock communities and spike-in controls. With a complex mock community spiked into an 146 environmental sample, uninfected control cells comprised as much as 16% of the total 147 community and four of the ten most abundant community members. Despite their dominance, 148 mock community control sequences were not observed in any fusion products (Fig. 3b). 149 Additionally, we ensured specificity by adding uninfected *E. coli* cells to replicate environmental 150 epicPCR reactions, to control for non-specific interactions. Although uninfected control 151 sequences were occasionally observed in the fusion data, these non-specific interactions were not 152 consistently associated with any specific phage sequence. Thus, non-specific associations were 153 removed by requiring that virus-host sequence pairs be observed in at least three libraries to be 154 considered a positive interaction (Fig 3c).

EpicPCR yielded 8,319 fusion amplicon sequences, containing 27 unique hosts (identical 168 rRNA gene sequence), 40 unique phages (identical RNR), and 95 unique phage-host interactions (Fig. 4, Table S3). RNR genes from epicPCR formed three distinct phylogenetic clades (Fig. 4a). These genes shared greatest homology with T7-like Chesapeake Bay cyanopodoviral isolates S-CBP1, S-CBP3, and S-CBP4 but had less than 80% nucleotide identity to any reference sequence (Table S2). Recruiting viral metagenome reads to the RNR amplicons

demonstrated these viral RNR sequences were not abundant in our samples, which highlights thesensitivity of this method (Extended Data Fig. 6).

163 Unexpectedly, these RNR sequences were most frequently associated with Actinobacteria 164 (Fig. 4) and are hereafter referred to as CSP-like Actinophage. Approximately 80% (31/40) of 165 the RNR sequences were linked to a single host (Actinobacteria member *Rhodoluna*), identifying 166 *Rhodoluna* as the primary host for these cyanopodoviral-like populations. We looked for other 167 cyanophage genes (e.g. *psbA*) on the same contig as these genes to investigate their genomic 168 context. However, RNR sequences failed to assemble into longer contigs, likely due to their low abundances and micro-diversity<sup>37</sup>. Other observations provide indirect evidence that these CSP-169 170 like Actinophage populations differ from their cyanophage counterparts. From May to December 171 2018, Chesapeake Bay *Rhodoluna* and CSP-like Actinophage abundances (Extended Data Fig. 172 3d) varied more similarly to each other (32 and 31-fold, respectively) than to cyanopodoviruses 173 (8-fold change) and Cyanobacteria (15-fold change). This suggests CSP-like Actinophage 174 abundances are associated with *Rhodoluna*, rather than with Cyanobacteria. After screening 175 single-amplified genome (SAG) libraries, no CSP-like Actinophage RNR gene sequences could 176 be amplified from ~300 Cyanobacteria SAGs. In contrast, one CSP-like Actinophage RNR gene 177 sequence was associated with an Alphaproteobacteria SAG, while a second was associated with 178 an Acidimicrobiia (phylum Actinobacteria) SAG library (Fig. 4a), providing further support that 179 poorly characterized CSP-like Actinophage interact with Actinobacteria.

Actinobacteria were also commonly predicted hosts of Chesapeake Bay viral populations using *in silico* approaches. We applied a variety of bioinformatics approaches to identify hosts for observed viral populations. Markov model-based prediction resulted in significant (p < 0.05) host predictions for about half of the viral populations from either reference genomes or

184 metagenome assembled genomes (MAGs; Extended Data Fig. 7). Four of the ten most frequently 185 predicted hosts were Chesapeake Bay MAGs, all of which were classified as Actinobacteria. 186 Among Chesapeake Bay MAGs, Actinobacteria were the predicted hosts for more viral 187 populations (57%) than all other MAGs combined even though Actinobacteria only comprised 188 25% of MAG classifications overall (Extended Data Fig. 7, 8). 2,204 tRNA sequences were 189 found in the viral populations, but this resulted in only 40 perfect matches between viral 190 populations and MAGs or reference genomes to allow for host inference. A majority of MAGs 191 harbored at least one CRISPR spacer, although only three MAGs had CRISPR spacers matching 192 observed viral populations. A comparison between bioinformatics methods demonstrates that 193 each approach identifies different and largely non-overlapping interactions within the vast virus-194 host interaction space (Extended Data Fig. 9). However, Actinobacteria were consistently 195 observed as putative hosts of viral populations in the Chesapeake Bay with all methods except 196 CRISPR analysis, suggesting significant viral pressure on Actinobacteria populations in this 197 environment throughout the year.

198 We also investigated Chesapeake Bay microbial susceptibility to viruses in vitro by 199 comparing growth of bacterial populations in incubations of water samples with and without 200 active viruses (Extended Data Fig. 10). Alphaproteobacteria were the most frequently identified 201 susceptible bacterial populations, and Flavobacteria displayed an interesting split between the 202 susceptible Flavobacteriaceae and the resistant Cryomorphaceae (Extended Data Fig. 10c). 203 Actinobacteria lineages Acidimicrobiia and Actinobacteria (class level) were also among the 204 most commonly identified susceptible taxa (Extended Data Fig. 10c). Three *Rhodoluna* spp. 205 were among these susceptible populations, including the primary *Rhodoluna* host identified by 206 epicPCR and a population that differed from it by one nucleotide (Extended Data Fig. 10c).

Previously, Actinobacteria was the most frequently identified host in an analysis of 2,000
freshwater phage genomes.<sup>38</sup> Although these phages lacked homology to the Chesapeake Bay
amplicons, this suggests that Actinobacteria may be under substantial viral pressure across both
freshwater and estuarine environments. More work is needed to understand how viruses shape
Actinobacteria abundance, diversity, evolution and metabolic processes.

212

## 213 In situ viral-host interaction dynamics reveal ecologically differentiated viral clades

214 EpicPCR results suggest RNR viral clades are ecologically differentiated. Interactions 215 between Clade I CSP-like Actinophages and their hosts peaked in May before largely 216 disappearing during the summer (Fig. 4B). Clade I phage sequences were most frequently 217 associated with multiple hosts (4/6 sequences) and the greatest number of hosts (Fig. 4A, C). 218 Host ranges for Clade I phages were significantly correlated with CSP-like Actinophage 219 abundances throughout 2018 ( $r_s = 0.69$ ; p = 0.002). Thus, broader host ranges of Clade I phages 220 could be a byproduct of increased contact rates resulting from higher viral abundances. Broader 221 host ranges may also provide an advantage for late spring viral populations under dynamic 222 conditions characterized by large freshwater influxes.

In contrast, Clade II and III phage-host interactions were primarily observed during the summer (Fig. 4B). Most (7/8) Clade II phage sequences were associated with a single host, while clade III sequences were split between single (11/19) and multiple (8/19) hosts. Neither clade displayed broad host ranges or correlations between host range and viral abundance like the Clade I phage sequences. Interactions between *Rhodoluna* and Clade II and Clade III phage populations alternated at most time points (Fig. 5), which may be a reflection of tidal influences on phage and/or host diversity. Supporting this hypothesis, Clade III interactions were negatively

correlated with fluorescent dissolved organic matter ( $r_s = 0.7$ ; p = 0.003), influenced by tide at our sample site<sup>39</sup>, and positively correlated with salinity ( $r_s = 0.6$ ; p = 0.02) (Table S4). Although mean tidal amplitude in the Rhode River is 0.3 m, water level is strongly influenced by weather conditions<sup>40</sup>, which might explain the lack of correlation with water level. This suggests Clade II and III phage populations may represent riverine and estuarine diversity, respectively. However, we cannot rule out other possible explanations, such as changes to host physiology with different environmental conditions<sup>41</sup> or antagonistic evolution between phage and host<sup>42</sup>.

237

#### 238 Host range associated with virus-host interaction lifespan and efficiency

Interactions captured by epicPCR illuminated differences in *in situ* interaction persistence and interaction efficiency between single-host (specialist) and multi-hosts (generalist) phage. Generalist phages were observed across a significantly greater number of sample dates  $(4.5 \pm 1.9)$ weeks) than specialist phages  $(2.5 \pm 1.0 \text{ weeks}; p = 0.002; \text{ Fig. 6A})$ . This suggests that broader host range increases the 'interaction lifespan' of phage, a possible evolutionary advantage in highly dynamic systems where viral infections appear to be largely ephemeral.<sup>43,44</sup>

245 EpicPCR and 16S rRNA marker gene analyses also revealed differences in the minimum 246 abundance of a host when it was associated with generalist or specialist phages. Virus-host 247 interactions could be detected by epicPCR at significantly lower host abundances for specialist 248 than generalist phages (p = 0.001; Fig. 6B). The minimum host abundance across time-points 249 where the host was associated with generalist CSP-like Actinophage ranged from 0 - 2.1% of the 250 community (0.63% mean minimum abundance). By comparison, the minimum abundance of 251 hosts associated with specialist CSP-like Actinophage ranged from 0 - 0.33% of the community 252 (0.11% mean minimum abundance). Thus, epicPCR was capable of detecting virus-host

253 interactions even when the host went undetected in the 16S rRNA gene libraries. Five of the 254 eight hosts associated with specialist phages were also associated with generalist phages, 255 suggesting that host identity was not inherent to these observed threshold differences. Instead, 256 observed minimum host abundance differences for specialist and generalist phages could be 257 attributed to increased infection efficiencies among specialist phages compared to generalist 258 phages. Supporting these findings, previous work identified fitness costs associated with increased host range in cultivated phage-host systems<sup>45</sup>, and trade-offs in viral infection 259 efficiency with increasing host range<sup>11</sup>. Considering these differences, there appears to be a 260 261 trade-off between 'interaction lifespan' and infection efficiency for CSP-like Actinophages in the 262 Chesapeake Bay.

263

#### 264 Discussion

265 EpicPCR can complement existing methods in viral ecology by linking viral and bacterial 266 populations in situ. EpicPCR targets specific viral populations, requires little more than standard 267 molecular biological equipment, and is cost-effective, making it ideally suited to probe specific 268 virus-host interactions with high temporal and genetic resolution. This targeted approach is more 269 cost-effective than non-specific approaches like proximity ligation or single-cell genomics when 270 studying a specific viral target. EpicPCR can complement large-scale metagenomic sequencing efforts<sup>6,46</sup> through host identification for uncultivated viruses, a notable advantage over 271 cultivation-based approaches, like viral tagging<sup>47</sup>. This method does not require cultivation to 272 273 reveal virus-host associations. However, as the dominant host was closely related to cultivated 274 Actinobacteria species (i.e. Rhodoluna), future work is needed to demonstrate the potential of 275 this technique to reveal virus-host relationships in uncultivated lineages. Digital PCR and single-

276 cell genomics can currently only accommodate hundreds to thousands of reactions and require 277 costly and/or specialized equipment. In contrast, epicPCR efficiently screens hundreds of 278 thousands of reactions within an emulsion, can be performed using standard molecular biology 279 equipment, and efficiently identifies positive associations when the vast majority of reactions are 280 expected to be negative. Although primer bias remains an issue for epicPCR, any of the commonly applied viral marker genes could be used<sup>48</sup>. Primers could be developed to target 281 abundant viruses, such as vSAG 37-F6<sup>49,50</sup>, or verify host predictions, such as potential archaeal 282 283 viruses<sup>51</sup>. Observations from epicPCR could test underlying assumptions of viral-host infection networks predicted from metagenomic time series.<sup>52</sup> Unlike cultivation-based methods or *in* 284 285 silico predictions, epicPCR can capture virus-host interactions when and where they occur, 286 revealing relationships with ecological and environmental factors. While epicPCR captures close 287 physical contact, it does not confirm active infections or differentiate between lytic or lysogenic infections, similar to limitations of single-cell genomics<sup>25</sup>. Pairing epicPCR with viral marker 288 gene expression<sup>43</sup> could confirm active infections. Host range capabilities of epicPCR paired 289 with quantitative methods of viral abundance, such as qPCR<sup>53</sup>, ddPCR<sup>50</sup> or polonies<sup>54</sup>, could be 290 291 used to estimate the impact of viruses on different microbial taxa.

In this study, we developed primers targeting a well-studied group of cyanophages, which yielded unexpected host associations and interaction dynamics. Our analysis revealed that these CSP-like RNR populations commonly interact with Actinobacteria, which was not predicted from homology searches of phage associated with cultured Actinobacteria, such as within PhagesDB<sup>55</sup>. Associations between observed viral populations and Actinobacteria were robust across methodologies, with all but one technique predicting Actinobacteria as a common host for observed viral populations. In addition, some phages interacting with Actinobacteria were also

299 associated with hosts from other phyla. Although most viruses are believed to have narrow host ranges, host enrichment may select against broad-host range phages<sup>56</sup>, and cultivated viruses 300 capable of infecting hosts spanning taxa at the order level have been reported<sup>57</sup>. Additionally, 301 Paez-Espino et al. (2016)<sup>58</sup> found 1% of viruses from metagenomic libraries had predicted hosts 302 303 spanning phyla based on CRISPR spacer and tRNA homology, including phages associated with 304 Actinobacteria. Furthermore, cultivation-based host range experiments often screen against 305 related hosts. We most frequently observed generalist phages associated with *Rhodoluna* 306 (Actinobacteria) and Halieaceae (Gammaproteobacteria), which would be unlikely to be 307 screened together. These broad host ranges could represent artifacts of the epicPCR method even 308 though we used controls in all environmental samples to remove spurious interactions. However, 309 observations of multiple different phages associated with the same hosts (Rhodoluna and 310 Halieaceae) suggest this is unlikely.

311

#### 312 Conclusions

313 Viruses influence biogeochemistry through interactions with hosts mediating key 314 microbial processes. Using a combination of shotgun metagenomics, marker gene analysis, 315 dilution experiments and a unique gene fusion technique, we investigated factors that influence 316 virus-host interactions within the largest estuary in the United States. We found Actinobacteria 317 populations are under substantial viral pressure within this ecosystem and their interactions with 318 specific phage clades seem to vary with tidally-influenced environmental factors. Virus-host 319 interactions revealed by epicPCR suggest that single- and multi-host phages have different 320 'interaction lifespans' and host abundance characteristics. If similar results are found broadly 321 across diverse viral lineages, these characteristics could represent an important trade-off shaping

viral evolution and might warrant separate parameters for generalist and specialist viruses in
ecological models. This combination of bioinformatics and experimental approaches provided
high genetic and temporal resolution of viral interactions with one of the most abundant
heterotrophic bacterial populations within this ecosystem that can be widely applied across
aquatic ecosystems to gain insight into viral ecology.

327

## 328 Materials and Methods

329 Sample Collection and preservation

330 Surface water samples were collected from May 2017 and May to December 2018 from 331 the mouth of the Rhode River, a tidal estuary on the Western Shore of the Chesapeake Bay 332 (Edgewater, MD), off of the Smithsonian Environmental Research Center research pier (38.89 N 333 76.54 W). Samples were collected five times a day over three days from 05/16/17-05/18/17, 334 between 10:30 am and 4:30 pm. Samples were also collected at 12:00 pm weekly from 5/24/18 335 to 8/9/18, on 8/23/18, then again weekly from 12/6/18 to 12/27/18. Briefly, 25 mL of water 336 sample was combined with 25 mL of 50% (v/v) sterile glycerol. Glycerol samples were stored on dry ice for transport back to Baltimore, MD and subsequently stored at -80°C until processing. 337 338 Approximately 120 mL of water sample was also collected per time point and filtered through a 339 0.2 µm PES membrane filter (Millipore, Inc.) for bacterial shotgun metagenomic analyses. 340 Filters were stored on ice for transport back to Baltimore, MD and kept at -80°C until 341 processing. Viral filtrates ( $< 0.2 \mu m$  fraction) were collected and incubated with FeCl<sub>3</sub> as previously described <sup>59</sup> during transport back to Baltimore, MD. Viral filtrates were filtered 342 343 through a 0.2 µm PES membrane filter (Millipore, Inc.) post-incubation, and filters were stored 344 in the dark at 4°C until processing for shotgun sequencing. Water conditions were recorded from

the continuous water monitoring station located at Smithsonian Environmental Research Center
(http://nmnhmp.riocean.com, Table S5).

347

348 Bacterial and viral gene counts

349 Bacterial and viral abundances were estimated by quantitative PCR. To create a standard 350 curve of gene copy number, viral 'Cyano SP'-like RNR genes were amplified from 351 environmental samples using primers Cyano II F and Cyano II R (0.3 uM final concentration 352 each, Table S6). Amplicons were run on a 1.5% agarose gel and visualized with SYBR Safe 353 DNA gel stain (Invitrogen, 1x final concentration). RNR bands were cut out, gel purified (Zymo, 354 Inc.), and cloned into chemically competent *Escherichia coli* cells using the Zero Blunt PCR 355 Cloning Kit (Thermo Scientific) following the manufacturer's protocol. Cells were grown overnight on LB + kanamycin (50  $\mu$ g mL<sup>-1</sup>) plates at 37°C and colonies were picked for 356 357 subsequent testing. Picked colonies were tested for amplification of the RNR gene. One colony 358 was serially diluted and grown on plates to determine gene copy number per dilution. The serial 359 dilution series was used to create a standard curve for use in qPCR. All standards and 360 environmental samples were run in triplicate. Three microliters of sample were combined with 361 UltraPure molecular grade water (Thermo, Inc.), SsoAdvanced Universal SYBR Green 362 Supermix (1x final concentration, Bio-Rad Laboratories, Inc.), Cyano II F primer (0.3 µM final 363 concentration), and Cyano II R primer (0.3 µM final concentration) to a final volume of 25 µL 364 (see Table S6 for primer sequences and citations). Samples were amplified on a CFX96 Real-365 Time PCR Detection System (Bio-Rad Laboratories, Inc.) with the following conditions: 366 denaturing at 98°C for 10 minutes; 45 cycles of denaturing at 98°C for 10 seconds, annealing at 367 52°C for 30 seconds, and extension at 72°C for 45 seconds; and a final extension of 72°C for 5

minutes. The same dilution series was used to create a standard curve for the 16S rRNA gene.
16S rRNA gene counts were assessed as above with primers PE\_16S\_U515F and 16S\_1114R
(Table S6).

371

372 Shotgun metagenomic library preparation, sequencing, and processing

373 To investigate bacterial diversity, we employed short-read shotgun metagenomics to the 374 bacterial fraction (> 0.2  $\mu$ m). DNA was extracted from 0.2  $\mu$ m filters for shotgun sequencing 375 from water samples collected on the following dates: 05/17/17, 05/18/17, 05/19/17, 05/31/2018, 376 06/28/18, 08/02/18, and 12/6/18. Additionally, two positive controls were processed, a Zymo 377 positive control community (Zymo Research) and E. coli, and one negative control (water). DNA 378 extraction was performed with the DNeasy PowerWater kit (Qiagen) following the 379 manufacturer's protocol with the following amendment: 20  $\mu$ L of proteinase K was combined 380 with 1 mL of solution PW1 in the bead tube. The bead tube was incubated at 65°C for ten 381 minutes prior to bead beating. Libraries were prepared with the Nextera DNA Flex Library Prep 382 kit (Illumina, Inc.) following the manufacture's protocol and sequenced on an Illumina MiSeq (2 383 x 300 bp) at the Genetic Core Research Facility at Johns Hopkins University. Sequences were quality filtered and trimmed with trimmomatic  $(v. 0.38)^{60}$ , assembled 384 with metaSPAdes  $(v. 3.13.1)^{61}$ . The assembly was used to generate metagenome assembled 385 genomes (MAGs) through metaWRAP (v1.2)<sup>62</sup> with metabat<sup>63</sup> and maxbin2<sup>64</sup>. Binning\_refiner<sup>65</sup> 386 was used to create the final set of MAGs with at least 50% completeness and less than 10% 387 contamination, as determined by CheckM<sup>66</sup> were also run in metaWRAP. MAGs were classified 388 using GTDB-Tk  $(v1.1.0)^{67}$ . 389

390

## 391 *16S rRNA gene amplicon library preparation, sequencing, and processing*

392	16S rRNA genes were amplified from 2017 and 2018 Chesapeake Bay surface water
393	samples (see Table S1) in a 25 $\mu L$ PCR reaction with the following conditions: three microliters
394	of column-purified DNA were combined with UltraPure molecular grade water (Thermo, Inc.),
395	10X buffer (1x final concentration), dNTPs (0.1mM each final concentration), 16S forward
396	primer 27F (0.3 µM final concentration), 16S reverse primer PE_16S_V4_E786_R (0.3 µM final
397	concentration), bovine serum albumin (0.02 mg/mL final concentration), and Phusion High-
398	Fidelity DNA Polymerase (0.5U; New England BioLabs, Inc.; see Table S6 for primer sequences
399	and citations). PCR reactions were combined with 150 $\mu L$ of 4% UMIL EM90 oil (4% UMIL
400	EM90 oil, 0.05% TritonX-100 v/v in mineral oil; Universal Preserv-A-Chem, Inc.) and
401	emulsified by vortexing at max speed (~2,700 rpm) for one minute on a Vortex Genie 2
402	(MoBio). Emulsions were loaded as 50 $\mu$ L aliquots and amplified with the following conditions:
403	denaturation at 94°C for 3 minutes; 33 cycles of denaturation at 94°C for 10 seconds, annealing
404	at 54°C for 30 seconds, and extension at 72°C for 45 seconds; and a final extension of 72°C for 5
405	minutes (C1000, BioRad Labs., Inc.). Samples were immediately removed upon completion of
406	amplification and stored at -20°C until the emulsion was broken.

PCR oil emulsions were broken with isobutanol as previously described<sup>68</sup>. Briefly, PCR
aliquots were pooled in a 1.5mL microcentrifuge tube and combined with 100 uL of sterile 5M
NaCl solution and 1 mL of isobutanol. Samples were vortexed briefly to mix and centrifuged at
16,000 x g for 1 minute. The bottom aqueous layer was retained, and DNA was purified by spin
column purification (Zymo, Inc.). DNA was eluted in 20 uL of Tris-HCl and stored at -20°C.
Purified DNA was run on a 1.5% agarose gel (UltraPure Agarose, ThermoFisher

413 Scientific) and visualized with SYBR Safe DNA gel stain (Invitrogen, 1x final concentration).

414	The gel was run in 1X TBE buffer (Alfa Aesar) at 4 V/cm. 16S rRNA gene bands were
415	visualized under blue light excitation, extracted, and gel purified (Zymo, Inc.) Purified DNA was
416	eluted into 20 $\mu$ L of Tris-HCl and stored at -20°C until further processing.
417	Barcodes and Illumina adapters were added to 16S rRNA gene amplicon products in two
418	subsequent limited PCR steps. Barcodes were added as follows: two microliters of purified DNA
419	were combined with UltraPure molecular grade water (Thermo, Inc.), 10X buffer (1x final
420	concentration), dNTPs (0.1mM each final concentration), 16S forward primer
421	PE_16S_V4_U515F (0.3 $\mu$ M final concentration), 16S rRNA gene reverse primer with 8-mer
422	barcodes PE_IV_XXX (0.3 $\mu$ M final concentration), and Phusion High-Fidelity DNA
423	Polymerase (0.5U; New England BioLabs, Inc.; see Table S6 for primer sequences and
424	citations). Samples were amplified with the following conditions: denaturing at 98°C for 30
425	seconds; 8 cycles of denaturing at 98°C for 10 seconds, annealing at 54°C for 30 seconds, and
426	extension at 72°C for 45 seconds; and a final extension of 72°C for 5 minutes. DNA was purified
427	by spin column purification (Zymo, Inc.) and eluted into 20 $\mu$ L Tris-HCl. Illumina adapters were
428	then added as above with the following primers: Illumina adapter forward primer PE-III-PCR-F
429	(0.3 $\mu$ M final concentration) and Illumina adapter reverse primer Barcode_Rev (0.3 $\mu$ M final
430	concentration) (see Table S6 for primer sequences and citations). Samples were amplified with
431	the following conditions: denaturing at 98°C for 30 seconds; 5 cycles of denaturing at 98°C for
432	10 seconds, annealing at 54°C for 30 seconds, and extension at 72°C for 45 seconds; and a final
433	extension of 72°C for 5 minutes. DNA was purified by spin column purification (Zymo, Inc.)
434	and eluted into 20 µL Tris-HCl.
435	16S rRNA gene amplicon products were quantitated on a Qubit 3.0 fluorometer

436 (Invitrogen) and three nanograms of DNA pooled per sample for sequencing. 16S rRNA gene

amplicon libraries were sequenced on an Illumina MiSeq (2 x 300 bp) at the Genetic Core
Research Facility at Johns Hopkins University. Sequence reads were processed in QIIME2<sup>69</sup>
using the DADA2 de-noising pipeline with the following parameters: trim left = 23 bases,
truncate = 200 bases, minimum fold-change of parent over abundance for chimera detection =
10. Taxonomic assignment was performed in QIIME2 with the Greengenes<sup>70</sup> database.

442

## 443 Virome shotgun metagenomic library preparation, sequencing, and processing

444 To investigate viral diversity, we employed short-/long-read hybrid metagenomic approaches<sup>71</sup> to the viral fraction ( $< 0.2 \mu m$ ) that had been incubated with FeCl<sub>3</sub> and filtered as 445 446 described above (Sample collection and preservation section). Samples from the following dates 447 were used for virome short-read shotgun metagenomic analysis: 05/17/17, 05/18/17, 05/31/18, 448 08/02/18, 12/14/18, 12/20/18. The two samples collected on 12/14/18, 12/20/18 were also long-449 read sequenced. After FeCl<sub>3</sub> incubation and filtration onto 0.2 µm filters, viruses were resuspended from filters with an ascorbic acid buffer as previously described<sup>59</sup>. Following 450 resuspension, viral particles were purified by cesium chloride gradient centrifugation <sup>72</sup> and 451 452 DNA extracted with Wizard Prep Columns (Promega, Corp.). Viral metagenome short-read 453 libraries were prepared using the NexteraXT kit (Illumina, Inc.) following the manufacturer's 454 protocol. For samples with > 0.16 ng/ $\mu$ L, samples were amplified with undiluted amplicon target 455 mix (ATM) at 15 cycles; for those with 0.1-0.16 ng/ $\mu$ L, samples were amplified with a 1:5 456 dilution of ATM at 18 cycles; for <0.1 ng/ $\mu$ L, samples were amplified with a 1:10 dilution of 457 ATM at 20 cycles. Short reads for all virome libraries were sequenced on an Illumina NovaSeq 458 S4 with 75M (target) 2x150bp reads at the JP Sulzberger Genome Center (Columbia University, 459 New York, NY). Additionally, the long-read libraries from December 2018 were prepared using

phenol:chloroform extraction protocol (dx.doi.org/10.17504/protocols.io.6cbhasn) and libraries
were prepared as previously described <sup>71</sup> with modifications and sequenced with an Oxford
Nanopore MinION instrument on a FLO-MIN106D R9 version Spot-ON flowcell (Rev D) at the
Ohio State University according to the manufacturer's instructions.

Short reads were cleaned and quality-trimmed with bbduk<sup>73</sup>; adapters, sequencing 464 465 artifacts, and PhiX sequences were removed (ktrim=r; k=23 mink=11; hdist=1; hdist2=1). Reads 466 were then quality-trimmed from both ends to remove bases with low quality scores (qtrim=rl; 467 trimg=20). Reads shorter than 30 bp (minlength=30), with Ns (maxns=0), or with an average 468 quality below 20 (mag=20) were discarded. The cleaned reads from each sample were then independently assembled with metaSPAdes (v. 3.13.1)<sup>61</sup> using k-mer sizes: 21, 33, 55, 77, and 469 470 -meta parameter. Long-reads were basecalled with Guppy v.2.3.1 (Manufacturer's tool) and 471 individual barcoded sample libraries were demultiplexed with the 'barcoder' function of Guppy. Long-reads were quality controlled with NanoFilt<sup>74</sup>, in which reads were filtered by quality 472 473 score (Q-score  $\geq 10$ ) and minimum length ( $\geq 1kb$ ), and finally 'headcropped' by 50bp to ensure 474 no remaining barcode sequence remained. These cleaned long-reads were used in two separate assembly scenarios. First, long-reads were assembled with Flye<sup>75</sup> in metagenomic mode (--meta). 475 Error-correction of the Flye assemblies where performed with Pilon v.1.23<sup>76</sup>, using the 476 477 corresponding short-reads mapping information [read recruitment was performed with BWA v.0.7.17<sup>77</sup>] to detect and reduce basecalling errors. Second, hybrid assembly, using both long-478 and short-reads was performed with SPAdes (v.3.13.1)<sup>61</sup>. Long reads were assembled through 479 hybrid assembly (with metaSpades hybrid option) and Flye<sup>75</sup>. In order to predict viral contigs in 480 the assembled datasets, assemblies were run through VirSorter<sup>78</sup> (v. 1.0.5) in the -virome mode 481 482 after upgrading its database with an expanded profile HMM database of viral proteins, mainly

483	from the GOV2.0 dataset <sup>14</sup> . Contigs that were resolved as categories 1, 2, 4 and 5 were retained,
484	and filtered by length $\geq$ 5kb ( $\geq$ 1.5kb, if circular). DeepVirFinder <sup>79</sup> was another tool for rescuing
485	additional viral contigs. We considered high-confidence viral contigs from DeepVirFinder to be
486	those of scores of $\ge 0.9$ with a p $\le 0.05$ , and lengths $\ge 5$ kb. These two sets of predicted viral
487	contigs (from both long- and short-read assemblies) resulted in 9,392 viral populations
488	(approximately species level genotypes) dereplicated from 14,848 virome (< $0.2\mu m$ ) and
489	metagenome (> 0.2 $\mu$ m) viral contigs using "ClusterGenomes" <sup>80</sup> using 95% average nucleotide
490	identity over 80% coverage of the shorter contig length for all contigs $> 5$ kb. <sup>14,81-83</sup>
491	To calculate the coverage of these viral populations, clean reads were mapped to the
492	Chesapeake Bay viral population database with Bowtie2 <sup>84</sup> in the non-deterministic and
493	sensitive mode. The output bam files were parsed using BamM
494	(https://github.com/Ecogenomics/BamM) to only keep the reads that covered over 70% of the
495	viral contig length, with over 75% read alignment length. Pysam v0.8.5
496	( <u>https://github.com/pysam-developers/pysam</u> ) was then used to filter out reads with <95%
497	identity. Trimmed pileup coverage "tpmean" for each contig was calculated using BamM and
498	then they were adjusted for each sample by metagenome size. The same read mapping strategy
499	was employed for RNR sequences upon constructing the rank abundance curves in Extended
500	Data Fig. 5.
501	To study seasonal diversity of Chesapeake Bay viral populations, all viromes were
502	randomly subsampled to 15M reads using bbmap "reformat" <sup>73</sup> with default parameters. The
503	subsampled read libraries were assembled using metaSPAdes (v. 3.13.1) as above. The resulting
504	assemblies were then processed with VirSorter and DeepVirFinder as above. Viral contigs

505 extracted using the same cutoffs as mentioned before were grouped into viral populations if they

506 shared  $\geq$  95% nucleotide identity across  $\geq$  80% of the shorter contig length. Subsequently, the 507 subsampled reads were recruited to the viral populations' representatives with Bowtie2, and the 508 same read-mapping cutoffs discussed above were applied to calculate sequence-depth adjusted coverage. Taxonomic assignment of viral populations was performed using vConTACT2<sup>80,85</sup> 509 510 with default parameters. First, the full proteome of every viral population representative was predicted using Prodigal  $(v2.6.3)^{86}$ . The protein set was then combined with all the proteins from 511 512 the phage and archaeal viruses in the NCBI RefSeq v88 release. The combined set of proteins 513 was then used as input for vConTACT2 to compute protein similarity overlaps (i.e. protein 514 clusters), and subsequently refined into genus-level equivalent 'viral clusters' (VCs). Using this 515 method, viruses are classified at the genus level. The resulting cluster file were imported and visualized in Cytoscape  $3.7.2^{87}$ . 516

517

## 518 Seasonal dynamics and diversity analyses

519 Seasonal bacterial communities were assessed from 16S rRNA gene amplicon libraries. Hierarchical clustering was performed in QIIME<sup>88</sup> with Bray-Curtis dissimilarity. Libraries were 520 521 sub-sampled at the smallest library size (150,000 observations) with ten jackknifed replicates. Community structure was visualized using heatmap.2 in the gplots<sup>89</sup> package in R. 522 523 Seasonal Chesapeake Bay viral communities were characterized by mapping reads to assembled contigs using Bowtie2<sup>90</sup> with default sensitive conditions. To compare viral 524 populations with a deeply sequenced virome library collected in 2012 at the same location<sup>32</sup>, 525 526 2017 and 2018 Chesapeake Bay virome contigs > 5kb were combined with contigs > 5kb from 527 the 2012 virome that were previously assembled. These merged datasets were then dereplicated 528 as described above, resulting in 21,784 viral populations. The Chesapeake Bay 2017 and 2018

529	reads were then mapped to this larger pool of viral populations. Only reads that mapped to
530	contigs with $\ge 90\%$ identity over $\ge 90\%$ of the read were retained. Viral populations were
531	considered present in the sample if retained reads mapped to $\geq$ 75% of the contig. In total, 10,858
532	of the 21,784 viral populations from the combined 2012, 2017, and 2018 samples were identified
533	in the 2017 and 2018 samples. Trimmed pileup coverage values were used as proxies of
534	observation counts for each population. Hierarchical clustering was performed in QIIME with
535	Bray-Curtis dissimilarity. Libraries were sub-sampled at the smallest library size (300,000)
536	observations with ten jackknifed replicates. Community structure was visualized in R as
537	described above.
538	Shannon diversity was calculated in QIIME for libraries with corresponding virome
539	samples (5/17/2017, 5/31/2018, 8/2/2018, 12/14/2018, 12/20/2018). Values were recorded as the
540	average of ten jackknifed replicates at a sampling depth of 10,000 observations for both 16S and
541	viral populations.
542	The distribution of Chesapeake Bay viral populations in ocean viral metagenomes was
543	investigated by mapping viral metagenome reads from Global Ocean Virome 2.0 database
544	libraries <sup>14</sup> to the Chesapeake Bay $>$ 5kb contigs as described above for the Chesapeake Bay 2012
545	viral metagenome sample. Additional libraries from previously published freshwater and
546	estuarine environments were similarly queried <sup>91-93</sup> , in addition to a viral metagenome from the
547	Damariscotta River Estuary, ME, USA (BioProject Accession No. PRJNA357591).
548	Fold-change in specific lineages between May and December 2018 were determined with
549	read recruitment (viral populations), qPCR (RNR genes), or a combination of qPCR and 16S
550	rRNA gene libraries. The difference in cyanopodovirus abundance (8-fold change) was
551	determined through read recruitment to populations classified as cyanopodovirus. The difference

in Cyanobacteria (15-fold change) and *Rhodoluna* (32-fold change) abundance was determined
by using the relative abundances from 16S rRNA gene libraries and the absolute abundance of
16S rRNA genes with qPCR. The difference in SCP-like Actinophage abundance (31-fold
change) was determined with qPCR.

556

557 Statistical analysis

558 The Kruskal-Wallis test for non-normal distributions was used to determine significance 559 of viral seasonal dynamics (p = 0.0003, Fig. S2b) and the growth of individual bacterial 560 populations in *in situ* incubations in the presence and absence of active viruses (FDR p < 0.05, 561 Extended Data Fig. 10). Spearman's Rho for non-normal distributions was used to determine the 562 p-value of the linear regression between the diversity of RNR genes with overall viral diversity (linear regression  $r^2 = 0.95$ , p = 0, Extended Data Fig. 4), the correlation of Clade I phage with 563 564 qPCR measurements of CSP-like Actinophage abundances (Spearman's Rho = 0.69; p = 0.002), the correlation of Clade III with fluorescent dissolved organic matter (Spearman's Rho = 0.7; p = 565 566 0.003; Table S4) and salinity (Spearman's Rho = 0.6; p = 0.02; Table S4). The interaction 567 lifespan (p = 0.002; Fig. 6A) and minimum host abundance (p = 0.001; Fig. 6B) differences 568 between generalist and specialist phage were determined with a two-tailed T Test for normally 569 distributed data. The calculation of minimum abundance of hosts associated with generalist 570 phage ranged from 0 - 2.1% of the community (0.63% mean minimum abundance, n = 21), 571 while the calculation of minimum abundance of hosts associated with specialists ranged from 0 -572 0.33% of the community (0.11\% mean minimum abundance, n = 8). The host was undetectable 573 in 5 out of 29 cases, but omitting these observations had no effect on the analysis.

## 575 Ribonucleotide reductase homology analysis

576 Open reading frames (ORFs) were called for all Chesapeake Bay > 5kb contigs using MetaGeneMark<sup>94</sup>. RNR alpha subunit genes were identified by BLASTx query of ORFs (e-value 577 < 1E-10) against a database of Uniref50<sup>95</sup> RNR cluster representative sequences. Putative RNR 578 genes were translated to amino acid sequences and aligned with MAFFT (v. 7.453)<sup>96</sup>. Aligned 579 sequences were visually inspected in Geneious<sup>97</sup> v. 9.1.5 for the presence of conserved catalytic 580 581 residues C439, E441, C462 indicative of RNR proteins (amino acid positions of Escherichia coli 582 nrdA gene product). Only sequences containing these conserved amino acids and spanning the 583 C462 to P621were retained. All retained sequences were queried against viral sequences in the 584 NCBI nr database by BLASTn to identify putative viral and host taxonomy. Top hits with a 585 percent identity  $\geq$  65% across at least 90% of the query sequence were recorded. All sequences 586 with top hits below this threshold were reported as 'Unclassified'. 587 Chesapeake Bay RNR amplicon sequences identified as CSP-like Actinophage by 588 epicPCR were queried against 2,034 freshwater phage genomes by BLASTn with an e-value 589 cutoff of 0.001. None of the freshwater phage genomes shared homology with the RNR 590 amplicon sequences.

591

## 592 Bioinformatics host prediction

593 Markov-model based predictions with WIsH<sup>98</sup> was used to determine potential hosts for 594 viral contigs greater than 10 KB in length. The entire dataset used in the WIsH analysis included 595 viral contigs and metagenome assembled genomes (MAGs >50% completion, <10% 596 contamination) from the Rhode River, along with the viral and hosts reference genomes from the 597 WIsH benchmark dataset, downloaded from NCBI. Because an appropriate null model is not

598 available for this unique estuarine environment, we assume that for every bacterial model, the set 599 of phage genomes in the dataset for which it is a host is negligible compared to the set of phage 600 genomes that for which it is not a host. Reference host and viral genome relationships in the 601 benchmark dataset were used to verify the performance under these assumptions. Running the 602 benchmark dataset with the null model yielded similar results as those previously reported (58%, 603 as compared to 63%, accuracy reported at the genus level). Thus, we applied the same 604 assumption for the null model to the analysis of environmental data. The analysis required top 605 hits to have a likelihood value in the top 5% of all calculated likelihoods ( $p \le 0.05$ ). Of the 606 9,392 viral contigs from the Chesapeake Bay, almost half had top predictions that were within 607 the top 5% of all tested likelihoods. Of significant hits, 80% of top predictions of hosts for 608 environmental viral population mapped to reference genomes and 20% of significant top 609 predicted hosts were MAGs.

610 Hosts for viral populations were determined by transfer RNA (tRNA) homology search. 611 tRNA genes were predicted from the 2017 and 2018 viral populations using ARAGORN 612  $(v1.2.38)^{99}$ . Matching tRNA sequences in either the MAGs or GenBank (nt database downloaded 613 on Nov. 6, 2020) was determined by BLAST<sup>100</sup> with 100% identity and coverage. The host for 614 each viral population was taken from either a match to a viral sequence, in which case the host of 615 the matching virus was assumed as the host, or a matching bacterial sequence. NCBI taxonomy 616 was used to get genus and phylum level taxonomic information for each match.

Potentially interacting viral populations were also determined for MAGs through
 CRISPR spacer matches to the observed viral populations. CRISPRs were identified within each
 metagenome assembled genome with CRISPRFinder online<sup>101</sup> using default settings. BLAST<sup>100</sup>
 was used to identify putatively interacting viral populations from the viral populations with

621 100% identity and coverage. Self-self matches were removed when contigs were identical622 between datasets.

623

624 epicPCR of environmental samples

625 Glycerol samples were thawed on ice and one mL was added to three replicate 1.5 mL 626 microcentrifuge tubes per sample. One replicate was left unamended, while the other two were 627 spiked with E. coli to approximately 0.1% and 1% of the bacterial community, respectively, to 628 identify false-positive interactions. A fourth replicate with 5% E. coli was processed for seven of 629 the time points. Samples were centrifuged at 25,000 x g for 10 minutes and resuspended after 630 supernatant removal to reduce free viral particles. Thirty microliters of each sample was 631 combined with UltraPure molecular grade water (Thermo, Inc.), 10X buffer (1x final 632 concentration), dNTPs (0.1mM each final concentration), Cyano SP F primer (1.0 µM final 633 concentration), Cyano SP R 519R fusion primers (R1 and R2 combined, 0.01 µM each final 634 concentration), S-\*-Univ-1100-a-A-15 16S reverse primer (1.0 µM final concentration), bovine 635 serum albumin (0.02 mg/mL final concentration), Tween-20 (0.8% v/v final concentration), and 636 Phusion High-Fidelity DNA Polymerase (1.5U; New England BioLabs, Inc.) to a final volume of 637  $75 \,\mu\text{L}$  (see Table S6 for primer sequences and citations). PCR reactions were combined with 450 638 uL of 4% UMIL EM90 oil (4% UMIL EM90 oil, 0.05% TritonX-100 v/v in mineral oil; 639 Universal Preserv-A-Chem, Inc.) and emulsified by vortexing at max speed (~2,700 rpm) for one 640 minute. Emulsions were loaded as 50 µL aliquots and amplified with the following conditions: 641 denaturation at 94°C for 3 minutes; 33 cycles of denaturation at 94°C for 10 seconds, annealing 642 at 54°C for 30 seconds, and extension at 72°C for 45 seconds; and a final extension of 72°C for 5

643	minutes (C1000, BioRad Labs., Inc.). Samples were immediately removed upon completion of
644	amplification and stored at -20°C until the emulsion was broken as described above.

#### 646 Confirmation of blocking primer efficacy prior to nested PCR amplification

647 Blocking primers are used during nested PCR to prevent the annealing and amplification of unfused genes from the emulsion  $PCR^{26}$ . Prior to enriching our fusion products via nested 648 649 PCR (detailed below), we tested the ability of the blocking primers to prevent amplification of 650 unfused products at relevant 16S and RNR gene concentrations. To simulate our nested PCR 651 conditions, RNR and 16S rRNA gene copy numbers were quantified in all epicPCR reactions 652 post-cleanup by qPCR as described above. Next, RNR genes were amplified from Rhode River 653 samples with primers Cyano II F and Cyano II R 519R primers (0.3 µM final concentration 654 each, Table S6). 16S rRNA genes were similarly amplified with primers 27F and 1492R (0.3 µM 655 final concentration each, Table S6). Unfused RNR and 16S rRNA gene amplicons were 656 combined at copy numbers equal to the highest observed across all samples after epicPCR as 657 determined by qPCR. Finally, the unfused RNR and 16S rRNA gene mixture was run through 658 nested PCR and barcoding PCR reactions with all environmental samples as described below. 659 Samples were run on a 1.5% agarose gel and visualized with SYBR Safe DNA gel stain 660 (Invitrogen, 1x final concentration). Fusion products were observed in the unfused mix control 661 when no blocking primers were used. However, with the addition of blocking primers no fusion 662 products were detected in the unfused mix control. We therefore proceeded with nested PCR amplification of the Chesapeake Bay fusion products with blocking primers (see below). 663 664

#### 665 Enrichment of viral-host fused amplicons

666	Fused amplicons were enriched by nested PCR. Three microliters of column-purified
667	DNA were combined with UltraPure molecular grade water (Thermo, Inc.), 10X buffer (1x final
668	concentration), dNTPs (0.1mM each final concentration), Cyano_SP_Nested_FA primer (0.3 µM
669	final concentration), 16S reverse primer PE_16S_V4_E786_R (0.3 $\mu$ M final concentration),
670	forward and reverse blocking primers U519F-block10 and U519R-block10 (1.0 $\mu$ M final
671	concentration each; to ensure no amplification of unfused genes, see above), and Phusion High-
672	Fidelity DNA Polymerase (0.5U; New England BioLabs, Inc.) to a final volume of 25 $\mu$ L (see
673	Table S6 for primer sequences and citations). Samples were amplified with the following
674	conditions: denaturing at 94°C for 30 seconds; 30 cycles of denaturing at 94°C for 10 seconds,
675	annealing at 54°C for 30 seconds, and extension at 72°C for 45 seconds; and a final extension of
676	72°C for 5 minutes. PCR reactions were cleaned by spin column purification (Zymo, Inc.) and
677	DNA eluted in 20 uL of Tris-HCl. Samples were barcoded following enrichment and cleanup by
678	PCR as described above with Cyano_SP_Nested_FB primer (0.3 $\mu$ M final concentration, Table
679	S6), reverse primer PE-IV-PCR-XXX with 8-mer barcodes (0.3 $\mu$ M final concentration, Table
680	S6), and forward and reverse blocking primers U519F-block10 and U519R-block10 (1.0 $\mu M$
681	final concentration each, Table S6). Barcoded samples were run on a 1.5% agarose gel
682	(Invitrogen) and visualized with SYBR Safe DNA gel stain (Invitrogen, 1x final concentration).
683	Fusion bands were cut out and gel purified (Zymo, Inc.). Fusion products were quantitated on a
684	Qubit 3.0 fluorometer (Invitrogen) and ten nanograms of DNA pooled per sample for
685	sequencing.

687 Sequencing and quality filtering of viral-host fusion amplicons

688 Fused amplicon products were sequenced on a PacBio Sequel with Sequel v3 chemistry 689 (University of Maryland Institute for Genome Sciences). Circular consensus sequences were 690 obtained from raw reads with the following parameters: minimum signal-to-noise ratio (SNR): 3, 691 minimum length: 500bp, minimum passes: 10, minimum read score: 0.75, minimum predicted 692 accuracy: 0.90. Consensus sequences passing these thresholds were oriented by searching for primer sequences using Cutadapt<sup>102</sup> with an error rate of 0.01. Only sequences passing this error 693 rate were retained. Following orientation, reads were demultiplexed in QIIME<sup>88</sup> with zero 694 695 allowed barcode mismatches. Fused gene products were split by gene and primer sequences identified with Cutadapt<sup>102</sup>. Only those sequences without mismatches in any primer sites (error 696 697 rate of 0.01) were retained for further analyses. Finally, amplicons were filtered by size to 698 remove truncated or chimeric reads. Only fusion sequences with RNR gene amplicons between 699 583 and 596bp and 16S rRNA gene amplicons between 249 and 262bp were retained. These 700 sequences were split into separate RNR and 16S genes and their primer sequences trimmed with Cutadapt<sup>102</sup>. 701

702

### 703 Analyzing viral and bacterial diversity from epicPCR fusion amplicons

Viral RNR and bacterial 16S rRNA gene sequences passing all quality filtering were clustered at 100% nucleotide identity with CD-HIT<sup>103</sup> to identify viral and bacterial populations. Only viral-host pairs observed in three or more libraries out of the 58 libraries total were considered positive interactions. 16S rRNA gene representative sequences were classified with the Ribosomal Database Project classifier Release 11.5<sup>104</sup>. RNR representative sequences were queried against the NCBI nr database via BLASTn to identify top hits. RNR gene sequences were aligned with MAFFT<sup>96</sup> using the L-INS-I setting. Sequences were trimmed to the region

711 G1380 to A2079 in the E. coli nrdA (RNR) gene and a maximum likelihood tree with 100 712 bootstrap replicates was made using the Jukes-Cantor substitution model in Phylogeny.fr (http://phylogeny.lirmm.fr)<sup>105</sup>. The three different clades (Fig. 4a) were characterized by five of 713 714 more sequences sharing >95% average nucleotide identity. Viral-host interaction matrices for 715 each week of the timeseries can be found in Supplemental File 1. Viral-host interaction networks 716 were visualized in Cytoscape. Chesapeake Bay RNR amplicon sequences identified as CSP-like 717 Actinophage by epicPCR were queried against 2,034 freshwater phage genomes by BLASTn 718 with an e-value cutoff of 0.001. None of the freshwater phage genomes shared significant 719 homology with the RNR amplicon sequences.

720

### 721 Viral dilution incubation experiments

722 One liter of surface water was collected in July 2019 from the mouth of the Rhode River 723 (Edgewater, MD) at the Smithsonian Environmental Research Center. 120mL was filtered 724 through two 0.2 µm PES membrane filter (Millipore, Inc.) and the viral filtrate was collected. 725 Half of the viral filtrate was autoclaved to eliminate active viruses. The remaining viral filtrate 726 was left unamended. Bacterial communities were resuspended off of the 0.2 um filters in 1mL 727 autoclaved viral filtrate by shaking for 10 minutes at medium speed on a MoBio Vortex Genie 2. 728 Resuspensions from each filter were pooled to create a single resuspended sample. Resuspended 729 cells were combined 1:120 with either autoclaved (no viruses) or un-autoclaved (with viruses) 730 viral filtrate. Samples were divided into twelve 20 mL replicates. Six replicates without viruses 731 and six replicates with viruses were incubated at room temperature. 500 µL aliquots were taken 732 at the beginning of the incubation and approximately every 24 hours over three days. 16S rRNA 733 and 'Cyano SP'-like RNR gene counts were quantified at each time point by qPCR as described

734 above. 16S rRNA gene amplicon libraries were created as described above for the initial 735 community and the community after 75 hours for all replicates. Growth of individual bacterial 736 populations over the course of the incubation was calculated as the fold-change in abundance 737 from T0 to T75 using qPCR results and relative abundances. OTUs that displayed growth in at 738 least four treatment replicates (without viruses or with viruses) were assessed for significantly 739 (FDR p < 0.05) greater growth in one of the treatments using a Kruskal-Wallis test. Putative 740 susceptibility was calculated for OTUs with significantly greater growth in one of the treatments 741 as (OTU Abundance Fold Change Without viruses)/(OTU Abundance Fold Change With 742 Viruses). Higher values indicate greater presumed susceptibility to viral-induced mortality.

743

### 744 Single-cell genomics

745 Samples from May 17, 2017 were sent to the Bigelow Single Cell Genomics Center (East 746 Boothbay, ME; https://scgc.bigelow.org/) for sorting of the cyanobacteria and prokaryotic 747 fraction into 384-well plates. Sorting was conducted with the red fluorescence as a function of 748 forward scatter, used to gate for "cyanobacteria", and Syto-9 stained DNA as a function of 749 forward scatter, used to gate for "all prokaryotes". Single cell genome amplification (SAG Generation 2 using WGA-X® amplification<sup>106</sup>) of sorted cells was conducted at the Bigelow labs 750 751 and returned to JHU for marker analysis. RNR and 16S rRNA marker genes were amplified from 752  $2 \mu L$  of a 1:100 dilution of the amplified genomes according to the protocol above, without the 753 use of emulsions. Positive amplicons were purified and sequenced with the forward primers at 754 the Genetic Resource Core Facility at JHU with an Applied Biosystems 3730xl DNA Analyzer. 755

756	Life Sciences Reporting Summary. Further information on experimental design and reagents is
757	available in the Life Sciences Reporting Summary.

759	Data availability.	Sequences associated	with 16S rRNA	libraries from	environmental	samples

and incubation experiments, bacterial and viral shotgun libraries, and fusion amplicons from

- repicPCR have been deposited in NCBI under BioProject accession number PRJNA599167.
- 762 Water physicochemical measurements and qPCR data have been deposited in the BCO-DMO

database under datasets 757405 and 821955. Datasets used in this analysis include Global Ocean

Virome 2.0 (GOV 2.0), NCBI non-redundant nucleotide database (nr), Tampa Bay metagenomic

- 765 libraries (BioProject Accession No. PRJNA28619, PRJNA47459, and PRJNA52403), and
- 766 Damariscotta River Estuary, ME, USA (BioProject Accession No. PRJNA357591).

767

768

# 770 **References**

- Suttle, C. A. Marine viruses—major players in the global ecosystem. *Nature Reviews Microbiology* 5, 801 (2007).
- Emerson, J. B. *et al.* Host-linked soil viral ecology along a permafrost thaw gradient.
   *Nature microbiology* 3, 870 (2018).
- Reyes, A., Semenkovich, N. P., Whiteson, K., Rohwer, F. & Gordon, J. I. Going viral:
  next-generation sequencing applied to phage populations in the human gut. *Nature Reviews Microbiology* 10, 607 (2012).
- Suttle, C. A. The significance of viruses to mortality in aquatic microbial communities.
   *Microbial ecology* 28, 237-243 (1994).
- Guidi, L. *et al.* Plankton networks driving carbon export in the oligotrophic ocean.
   *Nature* 532, 465-470, doi:10.1038/nature16942 (2016).
- Roux, S. *et al.* Ecogenomics and potential biogeochemical impacts of globally abundant
  ocean viruses. *Nature* 537, 689, doi:10.1038/nature19366 (2016).
- 784 7 Winget, D. M. *et al.* Repeating patterns of virioplankton production within an estuarine
   785 ecosystem. *Proceedings of the National Academy of Sciences* 108, 11506-11511 (2011).
- 786 8 Chen, X. W. *et al.* Tide driven microbial dynamics through virus-host interactions in the estuarine ecosystem. *Water Res* 160, 118-129, doi:10.1016/j.watres.2019.05.051 (2019).
- Flores, C. O., Meyer, J. R., Valverde, S., Farr, L. & Weitz, J. S. Statistical structure of
  host–phage interactions. *Proceedings of the National Academy of Sciences* 108, E288E297 (2011).
- Flores, C. O., Valverde, S. & Weitz, J. S. Multi-scale structure and geographic drivers of
   cross-infection within marine bacteria and phages. *The ISME journal* 7, 520-532 (2013).
- Jover, L. F., Cortez, M. H. & Weitz, J. S. Mechanisms of multi-strain coexistence in
  host–phage systems with nested infection networks. *Journal of theoretical biology* 332,
  65-77 (2013).
- Våge, S., Storesund, J. E. & Thingstad, T. F. Adding a cost of resistance description
  extends the ability of virus-host model to explain observed patterns in structure and
  function of pelagic microbial communities. *Environmental microbiology* 15, 1842-1852
  (2013).
- 800 13 Edwards, R. A., McNair, K., Faust, K., Raes, J. & Dutilh, B. E. Computational
  801 approaches to predict bacteriophage-host relationships. *Fems Microbiology Reviews* 40,
  802 258-272, doi:ARTN fuv04810.1093/femsre/fuv048 (2016).
- 80314Gregory, A. C. *et al.* Marine DNA viral macro- and microdiversity from pole to pole.804*Cell* **177**, 1109-+, doi:10.1016/j.cell.2019.03.040 (2019).
- Tully, B. J., Graham, E. D. & Heidelberg, J. F. The reconstruction of 2,631 draft
  metagenome-assembled genomes from the global oceans. *Scientific data* 5, 170203
  (2018).
- Burstein, D. *et al.* Major bacterial lineages are essentially devoid of CRISPR-Cas viral defence systems. *Nature communications* 7, 10613 (2016).
- 810 17 Hatfull, G. F. Dark matter of the biosphere: the amazing world of bacteriophage
  811 diversity. *Journal of virology* 89, 8107-8110 (2015).
- 812 18 Middelboe, M., Chan, A. M. & Bertelsen, S. K. in *Manual of Aquatic Viral Ecology*
- 813 (eds S. W. Wilhelm, M. G. Weinbauer, & C. A. Suttle) 118-133 (American Society of
  814 Limnology and Oceanography, 2010).

815	19	Deng, L. et al. Viral tagging reveals discrete populations in Synechococcus viral genome
816		sequence space. Nature 513, 242-+, doi:10.1038/nature13459 (2014).
817	20	Mosier-Boss, P. A. et al. Use of fluorescently labeled phage in the detection and
818		identification of bacterial species. Appl Spectrosc 57, 1138-1144, doi:Doi
819		10.1366/00037020360696008 (2003).
820	21	Allers, E. et al. Single-cell and population level viral infection dynamics revealed by
821		phage FISH, a method to visualize intracellular and free viruses. Environmental
822		<i>microbiology</i> <b>15</b> , 2306-2318 (2013).
823	22	Tadmor, A. D., Ottesen, E. A., Leadbetter, J. R. & Phillips, R. Probing individual
824		environmental bacteria for viruses by using microfluidic digital PCR. Science 333, 58-62
825		(2011).
826	23	Bickhart, D. M. et al. Assignment of virus and antimicrobial resistance genes to
827		microbial hosts in a complex microbial community by combined long-read assembly and
828		proximity ligation. Genome Biol 20, 18, doi:10.1186/s13059-019-1760-x (2019).
829	24	Stalder, T., Press, M. O., Sullivan, S., Liachko, I. & Top, E. M. Linking the resistome and
830		plasmidome to the microbiome. Isme Journal 13, 2437-2446, doi:10.1038/s41396-019-
831		0446-4 (2019).
832	25	Labonte, J. M. et al. Single-cell genomics-based analysis of virus-host interactions in
833		marine surface bacterioplankton. Isme Journal 9, 2386-2399, doi:10.1038/ismej.2015.48
834		(2015).
835	26	Spencer, S. J. et al. Massively parallel sequencing of single cells by epicPCR links
836		functional genes with phylogenetic markers. The ISME journal 10, 427 (2016).
837	27	Jang, H. B. et al. Taxonomic assignment of uncultivated prokaryotic virus genomes is
838		enabled by gene-sharing networks. <i>Nature biotechnology</i> <b>37</b> , 632-639 (2019).
839	28	Bench, S. R. et al. Metagenomic characterization of Chesapeake Bay virioplankton. Appl.
840		Environ. Microbiol. 73, 7629-7641 (2007).
841	29	Kan, J., Evans, S. E., Chen, F. & Suzuki, M. T. Novel estuarine bacterioplankton in
842		rRNA operon libraries from the Chesapeake Bay. Aquatic Microbial Ecology 51, 55-66
843		(2008).
844	30	Chen, F. et al. Diverse and dynamic populations of cyanobacterial podoviruses in the
845		Chesapeake Bay unveiled through DNA polymerase gene sequences. <i>Environmental</i>
846		<i>Microbiology</i> <b>11</b> , 2884-2892, doi:10.1111/j.1462-2920.2009.02033.x (2009).
847	31	Kan, J., Suzuki, M. T., Wang, K., Evans, S. E. & Chen, F. High temporal but low spatial
848		heterogeneity of bacterioplankton in the Chesapeake Bay. Appl. Environ. Microbiol. 73,
849		6776-6789 (2007).
850	32	Nasko, D. J. et al. Family A DNA polymerase phylogeny uncovers diversity and
851		replication gene organization in the virioplankton. Frontiers in microbiology 9, 3053
852		(2018).
853	33	Sakowski, E. G. <i>et al.</i> Ribonucleotide reductases reveal novel viral diversity and predict
854		biological and ecological features of unknown marine viruses. <i>Proceedings of the</i>
855		National Academy of Sciences 111, 15786-15791 (2014).
856	34	Dwivedi, B., Xue, B., Lundin, D., Edwards, R. A. & Breitbart, M. A bioinformatic
857		analysis of ribonucleotide reductase genes in phage genomes and metagenomes. <i>BMC</i>
858		Evolutionary Biology 13, 33, doi:10.1186/1471-2148-13-33 (2013).

- 859 35 Harrison, A. O., Moore, R. M., Polson, S. W. & Wommack, K. E. Reannotation of the
  ribonucleotide reductase in a cyanophage reveals life history strategies within the
  virioplankton. *Frontiers in microbiology* 10, 134 (2019).
- Suzek, B. E., Huang, H. Z., McGarvey, P., Mazumder, R. & Wu, C. H. UniRef:
  comprehensive and non-redundant UniProt reference clusters. *Bioinformatics* 23, 12821288, doi:10.1093/bioinformatics/btm098 (2007).
- 86537Martinez-Hernandez, F. *et al.* Single-virus genomics reveals hidden cosmopolitan and<br/>abundant viruses. *Nature communications* 8, 15892 (2017).
- Kavagutti, V. S., Andrei, A. S., Mehrshad, M., Salcher, M. M. & Ghai, R. Phage-centric
  ecological interactions in aquatic ecosystems revealed through ultra-deep metagenomics. *Microbiome* 7, 135, doi:10.1186/s40168-019-0752-0 (2019).
- Tzortziou, M. *et al.* Tidal marshes as a source of optically and chemically distinctive
  colored dissolved organic matter in the Chesapeake Bay. *Limnology and Oceanography* **53**, 148-159, doi:DOI 10.4319/lo.2008.53.1.0148 (2008).
- 40 Jordan, T. E., Pierce, J. W. & Correll, D. L. FLUX OF PARTICULATE MATTER IN
  874 THE TIDAL MARSHES AND SUBTIDAL SHALLOWS OF THE RHODE RIVER
  875 ESTUARY. *Estuaries* 9, 310-319, doi:10.2307/1351410 (1986).
- Kich Medium. *Applied and Environmental Microbiology* 45, 1316-1323, doi:Doi
  10.1128/Aem.45.4.1316-1323.1983 (1983).
- 42 Martiny, J. B., Riemann, L., Marston, M. F. & Middelboe, M. Antagonistic coevolution
  of marine planktonic viruses and their hosts. (2014).
- Sieradzki, E. T., Ignacio-Espinoza, J. C., Needham, D. M., Fichot, E. B. & Fuhrman, J.
  A. Dynamic marine viral infections and major contribution to photosynthetic processes
  shown by spatiotemporal picoplankton metatranscriptomes. *Nat Commun* 10, 1169,
  doi:10.1038/s41467-019-09106-z (2019).
- Moniruzzaman, M. *et al.* Virus-host relationships of marine single-celled eukaryotes
  resolved from metatranscriptomics. *Nat Commun* 8, 16054, doi:10.1038/ncomms16054
  (2017).
- Buffy, S., Turner, P. E. & Burch, C. L. Pleiotropic costs of niche expansion in the RNA
  bacteriophage Φ6. *Genetics* 172, 751-757 (2006).
- Sunagawa, S. *et al.* Structure and function of the global ocean microbiome. *Science* 348, 1261359, doi:10.1126/science.1261359 (2015).
- Borney, L. *et al.* Contrasting life strategies of viruses that infect photo-and heterotrophic
  bacteria, as revealed by viral tagging. *MBio* 3, e00373-00312 (2012).
- Adriaenssens, E. M. & Cowan, D. A. Using Signature Genes as Tools To Assess
  Environmental Viral Ecology and Diversity. *Applied and Environmental Microbiology*896
  80, 4470-4480, doi:10.1128/Aem.00878-14 (2014).
- Martinez-Hernandez, F. *et al.* Single-virus genomics reveals hidden cosmopolitan and abundant viruses. *Nat Commun* 8, 15892, doi:10.1038/ncomms15892 (2017).
- 899 50 Martinez-Hernandez, F. *et al.* Droplet Digital PCR for Estimating Absolute Abundances
  900 of Widespread Pelagibacter Viruses. *Front Microbiol* 10, 1226,
  901 10, 2200/5, i. 1, 2010, 01226 (2010)
- 901 doi:10.3389/fmicb.2019.01226 (2019).
- 90251Vik, D. R. *et al.* Putative archaeal viruses from the mesopelagic ocean. *Peerj* 5, e3428,903doi:10.7717/peerj.3428 (2017).

904 52 Jover, L. F., Romberg, J. & Weitz, J. S. Inferring phage–bacteria infection networks from 905 time-series data. Royal Society Open Science 3, 160654 (2016). 906 Brankatschk, R., Bodenhausen, N., Zeyer, J. & Burgmann, H. Simple Absolute 53 907 Quantification Method Correcting for Quantitative PCR Efficiency Variations for 908 Microbial Community Samples. Applied and Environmental Microbiology 78, 4481-909 4489, doi:10.1128/Aem.07878-11 (2012). 910 Baran, N., Goldin, S., Maidanik, I. & Lindell, D. Quantification of diverse virus 54 911 populations in the environment using the polony method. *Nat Microbiol* **3**, 912 doi:10.1038/s41564-017-0045-y (2018). 913 55 Russell, D. A. & Hatfull, G. F. PhagesDB: the actinobacteriophage database. 914 Bioinformatics 33, 784-786, doi:10.1093/bioinformatics/btw711 (2017). 915 Jensen, E. C. et al. Prevalence of broad-host-range lytic bacteriophages of Sphaerotilus 56 916 natans, Escherichia coli, and Pseudomonas aeruginosa. Applied and Environmental 917 Microbiology 64, 575-580 (1998). 918 57 Peters, D. L., Lynch, K. H., Stothard, P. & Dennis, J. J. The isolation and characterization 919 of two Stenotrophomonas maltophilia bacteriophages capable of cross-taxonomic order 920 infectivity. Bmc Genomics 16, 664, doi:10.1186/s12864-015-1848-v (2015). 921 Paez-Espino, D. et al. Uncovering Earth's virome. Nature 536, 425-+, 58 922 doi:10.1038/nature19094 (2016). 923 John, S. G. et al. A simple and efficient method for concentration of ocean viruses by 59 924 chemical flocculation. *Environmental microbiology reports* **3**, 195-202 (2011). 925 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina 60 926 sequence data. *Bioinformatics* **30**, 2114-2120 (2014). 927 Bankevich, A. et al. SPAdes: A New Genome Assembly Algorithm and Its Applications 61 928 to Single-Cell Sequencing. J Comput Biol 19, 455-477, doi:10.1089/cmb.2012.0021 929 (2012).930 Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP-a flexible pipeline for genome-62 931 resolved metagenomic data analysis. *Microbiome* 6, 158 (2018). 932 Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately 63 933 reconstructing single genomes from complex microbial communities. PeerJ 3, e1165 934 (2015).935 64 Wu, Y.-W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning 936 algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 937 605-607 (2015). 938 65 Song, W.-Z. & Thomas, T. Binning refiner: improving genome bins through the 939 combination of different binning programs. Bioinformatics 33, 1873-1875 (2017). Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: 940 66 941 assessing the quality of microbial genomes recovered from isolates, single cells, and 942 metagenomes. Genome research 25, 1043-1055 (2015). 943 Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to 67 944 classify genomes with the Genome Taxonomy Database. Bioinformatics 36, 1925-1927, 945 doi:10.1093/bioinformatics/btz848 (2019). 946 68 Schütze, T. et al. A streamlined protocol for emulsion polymerase chain reaction and 947 subsequent purification. Analytical biochemistry 410, 155-157 (2011). 948 69 Bolyen, E. et al. QIIME 2: Reproducible, interactive, scalable, and extensible 949 microbiome data science. Report No. 2167-9843, (PeerJ Preprints, 2018).

0.50	70	
950	/0	Desantis, I. Z. <i>et al.</i> Greengenes, a chimera-checked 165 rKNA gene database and
951	- 1	workbench compatible with ARB. Appl. Environ. Microbiol. 12, 5069-50/2 (2006).
952	/1	Warwick-Dugdale, J. <i>et al.</i> Long-read viral metagenomics captures abundant and
953		microdiverse viral populations and their niche-defining genomic islands. <i>Peerj</i> 7, e6800,
954		doi:10.7717/peerj.6800 (2019).
955	72	Hurwitz, B. L., Deng, L., Poulos, B. T. & Sullivan, M. B. Evaluation of methods to
956		concentrate and purify ocean virus communities through comparative, replicated
957		metagenomics. Environmental Microbiology 15, 1428-1440 (2013).
958	73	Bushnell, B. BBMap: A Fast, Accurate, Splice-Aware Aligner,
959		< <u>https://www.osti.gov/servlets/purl/1241166</u> > (2014).
960	74	De Coster, W., D'Hert, S., Schultz, D. T., Cruts, M. & Van Broeckhoven, C. NanoPack:
961		visualizing and processing long-read sequencing data. <i>Bioinformatics</i> 34, 2666-2669,
962		doi:10.1093/bioinformatics/bty149 (2018).
963	75	Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads
964		using repeat graphs. <i>Nat Biotechnol</i> <b>37</b> , 540-+, doi:10.1038/s41587-019-0072-8 (2019).
965	76	Walker, B. J. et al. Pilon: An Integrated Tool for Comprehensive Microbial Variant
966		Detection and Genome Assembly Improvement <i>Plos One</i> <b>9</b> e112963
967		doi:10.1371/iournal.pone.0112963 (2014)
968	77	Li H & Durbin R Fast and accurate short read alignment with Burrows-Wheeler
969	,,	transform <i>Bioinformatics</i> <b>25</b> 1754-1760 doi:10.1093/bioinformatics/btn324 (2009)
970	78	Roux S Enault F Hurwitz B I & Sullivan M B VirSorter mining viral signal from
971	10	microbial genomic data <i>Poori</i> <b>3</b> e985 doi:10.7717/neeri 985 (2015)
972	79	Ren L et al. Identifying viruses from metagenomic data using deen learning
973	17	Quantitative Riology doi:10.1007/s40484-019-0187-4 (2020)
974	80	Bolduc B. Vouens-Clark K. Roux S. Hurwitz B. L. & Sullivan M. B. iVirus:
075	00	facilitating new insights in viral ecology with software and community data sets
976		imbedded in a cyberinfrastructure <i>Isma Journal</i> <b>11</b> , 7, 14, doi:10.1038/ismai.2016.80
970		(2017)
079	01	(2017). Drum I. D. et al. Dotterns and ecological drivers of econ viral communities. Science
970	01	<b>349</b> 1261408 doi:10.1126/goioneg.1261408 (2015)
979	0 <b>7</b>	<b>546</b> , 1201496, doi:10.1120/science.1201496 (2015).
980	82	videsmood herizontal sone transfer. Bus Consuming 17, 020, dai:10.1186/s12864.016
981		widespread nonzontal gene transfer. Bmc Genomics 17, 950, doi:10.1180/S12804-010- $229(-\pi/2016)$
982	07	3280-X (2010).
983	83	Koux, S. <i>et al.</i> Minimum information about an Uncultivated Virus Genome (MIUVIG).
984	0.4	Nat Biotechnol 37, 29-37, doi:10.1038/nbt.4306 (2019).
985	84	Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nat Methods
986	~ <i>-</i>	9, 357-U354, doi:10.1038/Nmeth.1923 (2012).
987	85	Jang, H. B. <i>et al.</i> Taxonomic assignment of uncultivated prokaryotic virus genomes is
988		enabled by gene-sharing networks. <i>Nat Biotechnol</i> <b>37</b> , 632-+, doi:10.1038/s41587-019-
989		0100-8 (2019).
990	86	Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site
991		Identification. Bmc Bioinformatics 11, 119, doi:10.1186/1471-2105-11-119 (2010).
992	87	Smoot, M. E., Ono, K., Ruscheinski, J., Wang, P. L. & Ideker, T. Cytoscape 2.8: new
993		features for data integration and network visualization. <i>Bioinformatics</i> 27, 431-432,
994		doi:10.1093/bioinformatics/btq675 (2011).

995	88	Caporaso, J. G. et al. QIIME allows analysis of high-throughput community sequencing
996		data. Nature methods 7, 335 (2010).
997	89	Warnes, G. R. et al. gplots: Various R programming tools for plotting data. (2015).
998	90	Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nature
999		<i>methods</i> <b>9</b> , 357 (2012).
1000	91	Jasna, V., Parvathi, A. & Dash, A. Genetic and functional diversity of double-stranded
1001		DNA viruses in a tropical monsoonal estuary, India. Sci Rep-Uk 8, 16036,
1002		doi:10.1038/s41598-018-34332-8 (2018).
1003	92	McDaniel, L. D., Rosario, K., Breitbart, M. & Paul, J. H. Comparative metagenomics:
1004		natural populations of induced prophages demonstrate highly unique, lower diversity
1005		viral sequences. Environmental Microbiology 16, 570-585, doi:Doi 10.1111/1462-
1006		2920.12184 (2014).
1007	93	Allen, L. Z. et al. The Baltic Sea virome: Diversity and transcriptional activity of DNA
1008		and RNA Viruses. <i>Msystems</i> <b>2</b> , e00125-00116, doi:10.1128/mSystems.00125-16 (2017).
1009	94	Zhu, W., Lomsadze, A. & Borodovsky, M. Ab initio gene identification in metagenomic
1010		sequences. Nucleic acids research 38, e132-e132 (2010).
1011	95	Suzek, B. E., Huang, H., McGarvey, P., Mazumder, R. & Wu, C. H. UniRef:
1012		comprehensive and non-redundant UniProt reference clusters. <i>Bioinformatics</i> 23, 1282-
1013		1288 (2007).
1014	96	Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
1015		improvements in performance and usability. <i>Molecular biology and evolution</i> <b>30</b> , 772-
1016		780 (2013).
1017	97	Kearse, M. et al. Geneious Basic: an integrated and extendable desktop software platform
1018		for the organization and analysis of sequence data. <i>Bioinformatics</i> <b>28</b> , 1647-1649 (2012).
1019	98	Galiez, C., Siebert, M., Enault, F., Vincent, J. & Söding, J. WIsH: who is the host?
1020		Predicting prokaryotic hosts from metagenomic phage contigs. <i>Bioinformatics</i> <b>33</b> , 3113-
1021		3114 (2017).
1022	99	Laslett, D. & Canback, B. ARAGORN, a program to detect tRNA genes and tmRNA
1023		genes in nucleotide sequences. Nucleic Acids Res. 32, 11-16, doi:10.1093/nar/gkh152
1024		(2004).
1025	100	Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
1026		alignment search tool. J Mol Biol <b>215</b> , 403-410, doi:Doi 10.1006/Jmbi.1990.9999 (1990).
1027	101	Grissa, L. Vergnaud, G. & Pourcel, C. CRISPRFinder: a web tool to identify clustered
1028		regularly interspaced short palindromic repeats. <i>Nucleic Acids Res.</i> <b>35</b> , W52-W57.
1029		doi:10.1093/nar/gkm360 (2007)
1030	102	Martin M Cutadapt removes adapter sequences from high-throughput sequencing reads
1031	10-	<i>EMBnet journal</i> <b>17</b> 10-12 (2011)
1032	103	Fu L Niu B Zhu Z Wu S & Li W CD-HIT accelerated for clustering the next-
1033	100	generation sequencing data <i>Riginformatics</i> <b>28</b> 3150-3152 (2012)
1034	104	Wang O Garrity G M Tiedie I M & Cole I R Naive Bayesian classifier for ranid
1035	101	assignment of rRNA sequences into the new bacterial taxonomy <i>Annlied and</i>
1036		Environmental Microbiology 73 5261-5267 doi: Doi 10 1128/Aem 00062-07 (2007)
1037	105	Dereeper A <i>et al</i> Phylogeny fr robust phylogenetic analysis for the non-specialist
1038	100	Nucleic acids research <b>36</b> W465-W469 (2008)
1000		

1041 1042 106 Stepanauskas, R. *et al.* Improved genome recovery and integrated cell-size analyses of individual uncultured microbial cells and viral particles. *Nat Commun* 8, 10, doi:10.1038/s41467-017-00128-z (2017).

1043

1044 Correspondences and requests for materials should be addressed to Sarah Preheim

1045 (sprehei1@jhu.edu) and Eric Sakowski (e.g.sakowski@msmary.edu).

1046

## 1047 Acknowledgements

1048The authors would like to thank the Smithsonian Environmental Research Center and

1049 Katrine Lohan for providing access to their facilities during sample collection. This work was

1050 supported by the National Science Foundation Biological Oceanography (awards #1820652,

1051 #1829831 and #1756314) and a Gordon and Betty Moore Foundation Investigator Award

1052 (#3790). Part of this project was conducted using computational resources at the Maryland

1053 Advanced Research Computing Center (MARCC) and the Ohio Supercomputer Center (OSC)

1054 for High Performance Computing.

1055

## 1056 Author Contributions

1057 EGS and SPP conceived of the work. EGS conducted all field-work, epicPCR analysis, and

1058 incubation experiments associated with this work. EGS, KAW, and SPP conducted the

1059 experimental and computational analysis for the bacterial metagenome and 16S rRNA gene

1060 libraries. EGS, FT, AAZ, and OZ conducted experimental and computational analysis for the

1061 viral metagenomic libraries. EGS, KAW, and SPP conducted the bioinformatic host prediction.

1062 EGS wrote the manuscript. EGS, AAZ, OZ, MBS, and SPP edited the manuscript.

1063

## **Competing interests**

- 1065 The authors declare no competing financial interests.

### 1067 Additional information

- 1068 Supplementary information is available for this paper.
- *Reprints and permissions information* is available at www.nature.com/reprints.
- *Correspondence and requests for materials* should be addressed to EGS and SPP.

## 1072 Figure legends





- 1080 population was considered detected in a non-Chesapeake Bay sample if the sample's reads
- 1081 covered at least 70% of the genomic length of the viral population. Due to differences in the
- 1082 sequencing depth between libraries, the number of shared viral populations and the sum of their
- 1083 coverage (i.e. total sequencing depth; see Materials and Methods) were adjusted for each library
- 1084 size. For the stations that have multiple samples, the maximum value is reported.



1086

Figure 2. Seasonal dynamics of Chesapeake Bay bacterial and viral communities. Bacterial and viral surface water communities were collected from the mouth of the Rhode River, a subestuary of the Chesapeake Bay, in May 2017 and again in May, August, and December 2018. a.)
Bray-Curtis hierarchical clustering of bacterial communities based on 16S rRNA gene sequences shows that microbial communities are most similar within the same season. Libraries were sub-

1092	sampled at 150,000 observations with ten jackknifed replicates. Individual sequence variants are
1093	colored gray (present) or white (absent). b.) Bray-Curtis hierarchical clustering of viral
1094	communities from viral populations > 5kb collected across different years (see Materials and
1095	Methods) shows that seasonal community similarity exceeds within-year community similarity.
1096	Libraries were sub-sampled at 300,000 observations. Individual viral populations are colored
1097	gray. For the deeply-sequenced 2012 sample, only viral populations that were observed in the
1098	shallower 2017-2018 samples (4,328 of 12,736 viral populations, 34% of the 2012 populations)
1099	were included for clarity. c.) Bacterial and viral diversity (Shannon's H' index) for paired 16S
1100	rRNA gene and viral metagenomes were significantly correlated ( $r_s = 0.9$ , $p = 0.04$ ). Paired
1101	samples were from May 2017 (indicated by *), May 2018, August 2018, and December 2018 (n
1102	= 5 biologically independent samples). Viral metagenomes were subsampled to the smallest
1103	metagenome (15 M reads) prior to assembly for all Shannon's H' diversity analyses. Diversity
1104	indices represent the mean of jackknifed replicates ( $n = 10$ subsamples of the same initial
1105	community) at a sampling depth of 10,000 observations. Error bars are SD.
1106	



#### Conceptual Diagaram for Investigating in situ Host-Virus Interactions with epicPCR



a,

Figure 3. EpicPCR identifies phage-host interactions in the environment without cultivation. a.) Overview of the experimental design for epicPCR, identifying phage-host interactions through single-cell isolation in emulsion droplets and virus and host marker gene fusion. Left: individual cells are isolated within emulsion droplets and the genome of the host (blue circle) and virus (red curved lines) serve as the template for the gene fusion reaction. Middle: fusion PCR joins and

amplifies viral and host marker genes from actively infected cells within emulsion droplets. RNR

1115 (ribonucleotide reductase; red) and 16S (16S rRNA gene; blue) are joined through an

1116 overlapping primer sequence (light blue). Right: fused amplicons are sequenced and analyzed to

1117 identify the network of viral-host interactions in the environment. b-c.) Control experiments to

- 1118 test epicPCR specificity. Specificity of the method was tested by spiking Chesapeake Bay water
- samples with a mock community of ten OTUs (cultivated cells, inactivated) or a single
- 1120 uninfected host (E. coli) prior to emulsification. b.) Proportion of fusion sequences belonging to
- 1121 uninfected mock (red) and Chesapeake Bay (blue) community members as a function of their

- 1122 community abundance. Uninfected mock community sequences were not found to be associated
- 1123 within any fusion products. c.) Proportion of fusion products containing E. coli when spiked into
- 1124 Chesapeake Bay samples at 0.1% (n = 17), 1% (n = 17), and 5% (n = 7) of the community. Post-
- 1125 filtered sequences are those interactions observed in a minimum of three libraries. Box and
- 1126 whisker markers represent the minimum, first quartile, median, third quartile, and maximum
- 1127 values.
- 1128
- 1129



Figure 4. Abundance, diversity, and ecology of CSP-like Actinophage virus-host interactions as determined by epicPCR in the Chesapeake Bay from May to December 2018. a.) Maximum likelihood tree with 100 bootstrap replicates of Chesapeake Bay CSP-like Actinophage RNR gene sequences (G1380 to A2079 in *E. coli nrdA*, See Materials and Methods). Nodes are colored by bootstrap support: Black > 75%; Grey 50 – 75%; no color < 50%. Only RNR gene sequences that were observed with the same host (100% nucleotide identity of 16S rRNA gene) in at least three epicPCR fusion amplicon libraries were included. Phage RNR gene sequences

1138 formed three main clades where five or more RNR sequences shared > 95% nucleotide identity. 1139 Phage RNR sequences were clustered at 100% nucleotide identity, resulting in 40 unique 1140 sequences shown on the tree. Host taxonomy identified by 16S rRNA gene homology for each 1141 phage RNR gene sequence is indicated. In some cases, closely related viruses were associated with different hosts. Scale bar represents nucleotide substitutions per site. b.) Mean percentage of 1142 1143 total identified interactions observed at each sample time point by phage RNR clade. Phage-host 1144 interactions were investigated by epicPCR weekly from May to August 2018 and again weekly 1145 in December 2018. Clade I phage RNR sequences were primarily observed in late spring, while 1146 RNR sequences from Clades II and III were observed throughout the summer. c.) Aggregate 1147 network of Chesapeake Bay CSP-like Actinophage-host interactions from May to December 1148 2018. Host taxonomy is colored as in panel A and is represented with rectangular shapes. Phage 1149 taxonomy is represented by the color map on the left and is represented by v-shapes. Lines 1150 represent virus-host interactions observed with epicPCR, colored according to viral host range. 1151 Most CSP-like Actinophage RNR sequences were associated with a single Actinobacteria host 1152 16S rRNA gene identified as Luna-1 member Rhodoluna. 1153



1154 1155 Figure 5. Phage interactions with Rhodoluna host observed in the Chesapeake Bay from May to 1156 August 2018 and again in December 2018 as revealed by epicPCR. In total, 31 of 40 phage 1157 populations identified by epicPCR interacted with this host. a.) Maximum likelihood tree with 1158 100 bootstrap replicates of Chesapeake Bay phage ribonucleotide reductase (RNR) genes 1159 (G1380 to A2079 in E. coli nrdA, See Materials and Methods). RNR gene sequences were 1160 clustered at 100% nucleotide identity prior to phylogenetic analysis. Phage RNR gene sequences 1161 formed three main clades where RNR sequences shared > 95% nucleotide identity. Scale bar 1162 represents nucleotide substitutions per site. b.) Proportion of total observed interactions

- 1163 identified by epicPCR each week between Rhodoluna host and each identified phage RNR gene
- 1164 sequence. Colored wedges link phage populations to their corresponding dynamics over time.
- 1165 Wedge colors correspond to phage clades: blue: Clade I; orange: Clade II; grey: Clade III; no
- 1166 color: no assigned clade. Only the most abundant phage interactions were depicted for
- 1167 simplicity. c.) Proportion of phage interactions by clade with Rhodoluna host over time
- 1168 identified by epicPCR. Blue: Clade I; orange: Clade II; grey: Clade III.
- 1169



1172 Figure 6. Biological trade-offs and ecological patterns related to viral and bacterial interactions 1173 revealed by epicPCR. a.) Total interaction persistence observed by epicPCR for CSP-like 1174 Actinophage RNR sequences between May and December 2018 in the Chesapeake Bay. RNR 1175 sequences associated with a single host in the epicPCR amplicons were designated specialist phages (n = 26), while RNR sequences associated with multiple hosts in the epicPCR amplicons 1176 were designated generalist phages (n = 14) (\*\* = two-tailed T Test p = 0.002 for normally 1177 1178 distributed data). Normal distributions were determined by the D'Agostino and Pearson test (p >1179 0.05). Box and whisker markers represent the minimum, first quartile, median, third quartile, and maximum values. The mean is indicated by "+". b.) Minimum host abundance (16S rRNA gene 1180 1181 relative abundance) of interactions with specialist phage RNR gene sequences (n = 8) and 1182 generalist phage RNR gene sequences (n = 21) observed in epicPCR fusion amplicons from 1183 samples collected between May and December 2018 in the Chesapeake Bay (\*\* = two-tailed T 1184 Test p = 0.001 for normally distributed data). Minimum host relative abundances indicate the 1185 lowest observed host abundance across timepoints where the host was associated by epicPCR 1186 with a generalist or specialist phage. All hosts identified by epicPCR were present in at least 1187 three 16S rRNA gene amplicon libraries. Data points where host abundance equals zero (n =

- 1188 5/29 observations) indicate an observed phage-host interaction in epicPCR amplicons but no host
- 1189 detected in the 16S library for that time point. Omitting instances where the host was undetected
- 1190 in 16S rRNA gene amplicon libraries or using the next lowest abundance observation had no
- 1191 impact on the result. Normal distributions were determined by the D'Agostino and Pearson test
- 1192 (p > 0.05). Box and whisker markers represent the minimum, first quartile, median, third
- 1193 quartile, and maximum values. The mean is indicated by "+".