**1  Dynamic marine viral infections and major contribution to photosynthetic processes shown**
**2  by regional and seasonal picoplankton metatranscriptomes**
3
4  Sieradzki Ella T.[1], Ignacio-Espinoza J. Cesar[1], Needham David M.[2], Fichot Erin B.[1] and
5  Fuhrman Jed A.[1,*]
6
7  Author affiliations: (1) University of Southern California, (2) Monterey Bay Aquarium Research
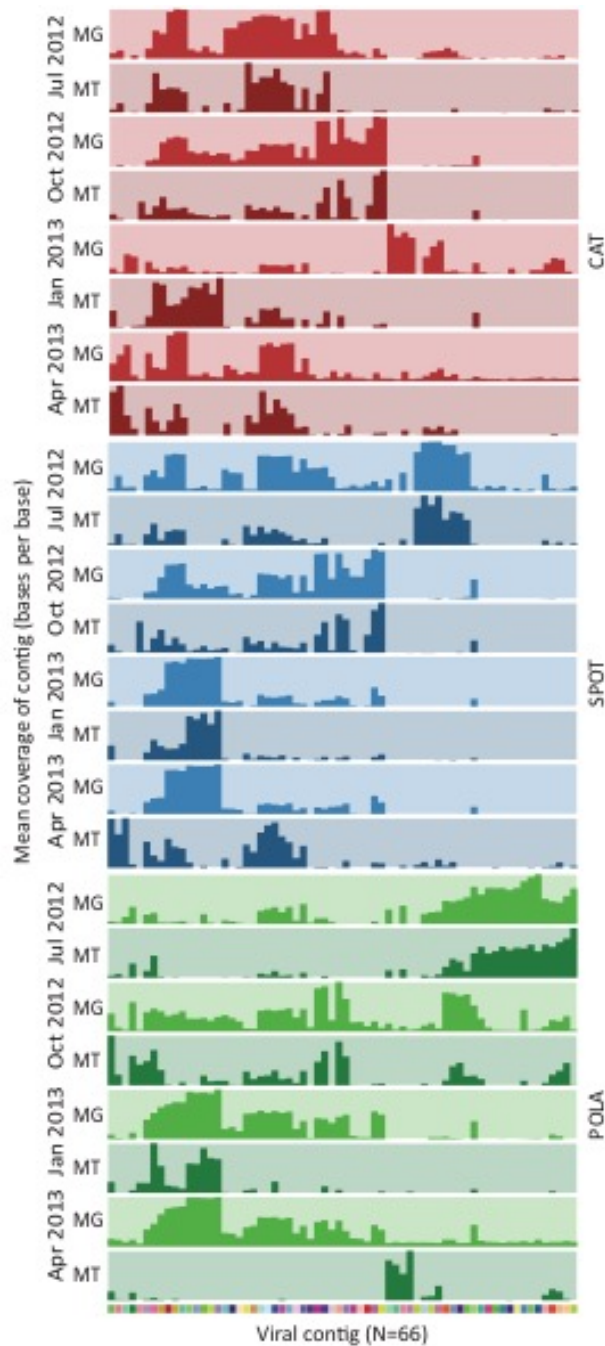8  Institute
9  * Corresponding author
10

**11  Viruses are an important top-down control on microbial communities, yet their direct**
**12  study in natural environments has been hindered by culture limitations[1-3]. The advance of**
**13  sequencing and bioinformatics over the last decade enabled the cultivation independent**
**14  study of viruses. Many studies focus on assembling new viral genomes[4-6] and studying viral**
**15  diversity using marker genes amplified from free viruses[7,8]. We used cellular**
**16  metatranscriptomics to study community-wide viral infections at three coastal California**
**17  sites throughout a year. Generation of and recruitment to viral contigs (> 5kbp, N=66)**
**18  allowed tracking of infection dynamics over time and space. Here we show that while these**
**19  assemblies represent viral populations, they are likely biased towards clonal or low**
**20  diversity assemblages. Furthermore, we demonstrate that published T4-like cyanophages**
**21  (N=50) and pelagiphages (N=4), having genomic continuity between close relatives, are**
**22  better tracked using marker genes. Additionally, we demonstrate determination of**
**23  potential hosts by matching infection dynamics with microbial community composition.**
**24  Finally, we quantify the relative contribution of various cyanobacteria and viruses to**
**25  photosystem-II *psbA* expression in our study sites. We show sometimes >50% of all**
**26  cyanobacterial+viral *psbA* expression we observed is of viral origin, which highlights the**
**27  proportion of infected cells and makes viruses a remarkable contributor to photosynthesis**
**28  and oxygen production.**
29

30  We sampled surface seawater in different seasons over three sites across the San Pedro Channel,
31  California, USA: The Port of Los Angeles (POLA), Santa Catalina Island Two Harbors (CAT)
32  and the San Pedro Ocean Time-series (SPOT). These sites represent a gradient of human impact
33  with POLA being the most impacted and SPOT resembling open ocean conditions. In all of these
34  sites free virus-like particles outnumber bacteria and archaea roughly 10:1 (sup. fig. 1). We
35  examined only the 0.2-1 µm size-fraction, which includes most bacteria, archaea and some
36  picoeukaryotes. Via assembly of metatranscriptomes, we obtained 1455 contigs longer than 5 kb
37  of which 57 (3.9%) were characterized as viral using virSorter and virFinder (see methods).
38  Additionally, a cross-assembly of the metatranscriptomic viral contigs with metagenomes of the
39  same samples (N=12) yielded 9 more contigs (mean length 26,563 bp) characterized as viral.
40  Most of the contigs represent dsDNA viruses (N= 65) as apparent from their presence in
41  metagenomes, but one appears to be an RNA virus possibly infecting a eukaryotic host. This
42  contig contained an RNA-dependent-RNA-polymerase whose nearest match in NCBI non-
43  redundant database was marine Antarctic phytoplankton RNA virus PAL_E4[9]. These 66 viral
44  contigs revealed varied patterns of presence (in metagenomes) and activity (in
45  metatranscriptomes) in the three sites over a year (fig. 1).

46    Active non-synchronized viral infection would manifest as recruitment to an entire contig in both
47    metagenome and metatranscriptomes of the same sample. We found that patterns of mean
48    coverage from metagenomes and metatranscriptomes of our assembled viral contigs usually
49    differed, not just between metagenomes and metatranscriptomes but also between dates and
50    locations, implying widespread boom-bust dynamics of infection. While some variation may be
51    due to synchronization known for some photosynthetic and heterotrophic bacteria in the
52    ocean[10,11] and for some of their phages[12], this explanation is less likely as samples were collected
53    from all sites within the same 4 hours morning-time window.
54    Some regional patterns were evident, e.g. some viral contigs were unique to the Port of LA (fig.
55    1), and that site always clustered separately from SPOT and CAT by Bray-Curtis similarity of
56    expression of viral contigs (sup. fig. 1B). This pattern corresponds to the difference in biotic
57    parameters between the port and the other sites (sup. fig. 2), though the port did not cluster
58    separately in microbial community composition by 16S-rRNA (sup. fig. 1A). The latter may
59    reflect offshore microbes brought in with the tide but less active than port organisms. Clustering
60    by metagenomic recruitment to viral contigs did not reveal consistent patterns by site or date
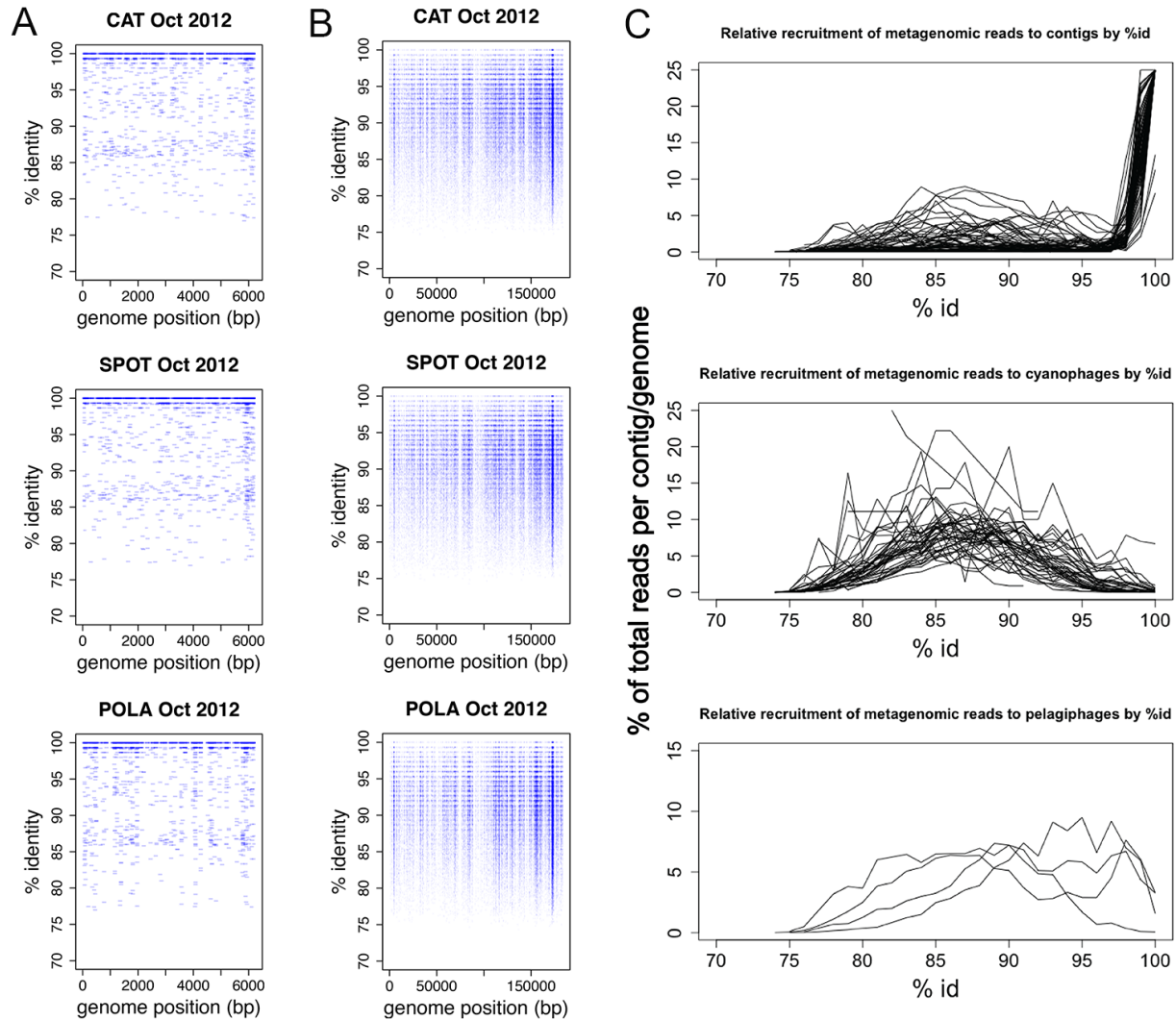61    (sup. fig. 1C).
62

63
64
65 **Figure 1**: Mean coverage of 66 viral contigs across three sites (Port of LA – POLA, San Pedro
66 Ocean Time-series – SPOT and Two harbors – CAT) and four dates (July 2012, October 2012,
67 January 2013 and April 2013) in metagenomes (MG) and metatranscriptomes (MT). The bar
68 heights are normalized to the highest mean coverage within the sample. Each cell in the color bar
69 on the bottom represents a contig and corresponds with the column above it in all samples. Mean
70 coverage was calculated excluding contig positions in the 4[th] quartile of coverage depth which
71 can be biased by recruitment localized to a small portion of the contig (sup. fig. 3).
72

73  Ephemeral infections dominated the assembled landscape, as 56 out of 66 of the contigs only
74  appeared in few metatranscriptomes, presumably reflecting sporadic infections. Persistent
75  infections (mean coverage >= 0.75x in at least 3 out of 4 samples per site, 10 out of 66) were
76  limited to CAT and SPOT except for one that was persistent in all three sites. Moniruzzaman et
77  al.[13] also recently demonstrated dominance of ephemeral dynamics in infections of marine
78  single-cell eukaryotes during an algal bloom. Bray-Curtis dissimilarity of the viral contigs within
79  each site was 80-100%, whereas the dissimilarity of microbial communities within site was
80  distributed around 50-70%. High dissimilarity indicates that even within site different viruses are
81  actively infecting in different seasons (sup. fig. 1D+E).
82  Moreover, assembled viral contigs appeared to be biased towards low-microdiversity (i.e. more
83  clonal) viruses. High diversity, extremely common in marine microorganisms[14], tends to break
84  assemblies created with either read-overlaps or DeBruijn graphs[15,16]. We expect that low virus
85  diversity could result from boom-bust lifestyle due to bottlenecks during "bursts". This might
86  lead to a method bias towards ephemerally infecting viruses. Indeed, all the viral contigs we
87  assembled in this study appear to have many nearly identical relatives but few moderately close
88  ones as shown by recruitment plots (most recruitment at 98-100% identity and little recruitment
89  at 90-97%, fig. 2C), while some of the published pelagiphages had recruitment along most of the
90  genome and high mean coverage at up to 100% identity and yet did not assemble (fig. 2C, sup.
91  table 1).
92  The recruitment plots also reveal a common pattern of recruitment to short fragments near 100%
93  identity whereas the rest of the genome or contig is only recruited to at lower percentage if at all
94  (sup. fig. 3). This pattern highlights two issues: (1) some genes are so conserved or so often
95  laterally transferred that their partial sequences cannot be used to identify which phage is present
96  and (2) that mean coverage of contigs could be highly biased by these conserved regions which
97  needs to be considered when evaluating abundance of the contigs and for coverage-based
98  binning of genomes.
99  A previous report indicated that Synechococcus phage genomes occur in discrete "clouds" with a
100 discontinuity in recruitment below ~95% identity[17]. While this pattern exists for some
101 cyanophage genomes, and we often saw some gaps in coverage at ~90-95% consistent with that
102 idea (sup. fig. 3), it is by no means the rule in our data, especially for pelagiphages (fig. 2C). We
103 also note that widely used recruitment algorithms only map reads with a local or end-to-end
104 match at a very high percent identity, and would therefore miss much genetic diversity that may
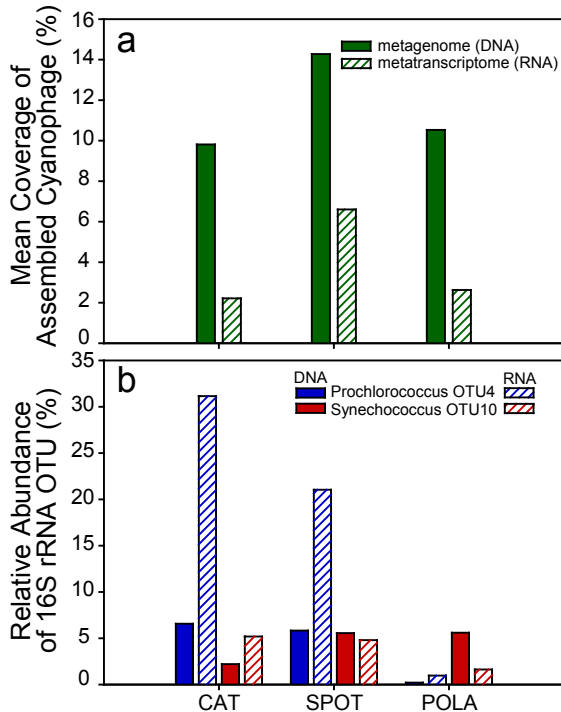105 be relevant (fig. 2B).

**Figure 2**: Metagenomic read recruitment to (**A**) an assembled cyanophage contig and (**B**) *Prochlorococcus* phage P-HM2 genome. Most recruitment to the assembled contig is at 99-100% identity (high density near 100% is not fully evident from the graph due to overlaps, see C), whereas P-HM2 reveals a genomic continuum. (**C**) Recruitment as a function of percent ID of reads demonstrates that assembled contigs mostly recruit at 100% ID and have few moderately close relatives (top) whereas published genomes of cyanophages reveal clouds of moderately close relatives but few matches near 100% (middle), and pelagiphages range from 100% down (bottom).

We were surprised not to find multiple cyanophage (especially myovirus) contigs, because such cyanophages belong to the family *Myoviridae*, some of the most common dsDNA viruses in the ocean[18] and we know this region has a diverse community of myoviruses and cyanobacteria[7,14]. Few of the assembled viral contigs contained myoviral marker genes (e.g. capsid protein gp23) (sup. Table 2). The only assembled contig that is with high certainty from a cyanophage is a putative podovirus (see below). Recruitment of reads to published cyanophage genomes revealed the likely reason for so few such contigs: high genomic diversity (fig. 2B) which probably broke

124    assemblies of T4-like cyanophages. We lacked assemblies despite persistent myovirus activity.
125    We assigned translated reads identified by a Gp23-HMM (Hidden Markov Model) to published
126    and assembled Gp23 proteins. Most versions of this marker gene from published genomes as
127    well as the nine assembled Gp23 ORFs were expressed persistently throughout all sites and dates
128    (sup. fig. 4). While the exact published genomes themselves were not present in our samples (fig.
129    2B), we posit that other T4-like cyanophages closely-related to those published are present and
130    persistently infecting their hosts.
131    Matching viral contigs and hosts is challenging, but we were able to use physiological
132    information and distributions among samples to make a likely match. Many cyanophages contain
133    a variety of genes that maintain photosynthetic activity in the host during infection, from "spare
134    parts" for photosynthetic reaction centers through regulation and optimization of those apparati[19].
135    In particular, viruses were shown to maintain photosystem II function during infection in order to
136    supply energy to the host, as transcription of host genes is shut down during infection and PS-II
137    proteins have a short lifetime[20,21]. Our assembled cyanophage contig contained genes coding for
138    photosystem-II protein D1 (*psbA*) and high-light induced protein (*hli*) reportedly widespread in
139    cyanophages[8]. The putative cyanophage from which this contig was derived was actively
140    transcribed (presumably infecting its host) in all three sites only in October 2012 (fig. 4A). The
141    cyanobacterial community by 16S-rRNA was dominated in October by two operational
142    taxonomic units (OTUs): one *Synechococcus* and one *Prochlorococcus*. Both OTUs were present
143    at SPOT and CAT in October, but only *Synechococcus* was also present at POLA (fig. 4B).
144    Thus, we propose that this assembled contig is from a phage that infects *Synechococcus* OTU 10
145    which has a 16S sequence over the amplified region 100% identical to *Synechococcus* CC9902
146    of clade IV. On a phylogenetic tree of PS-II D1, translated PS-II D1 of this phage clustered
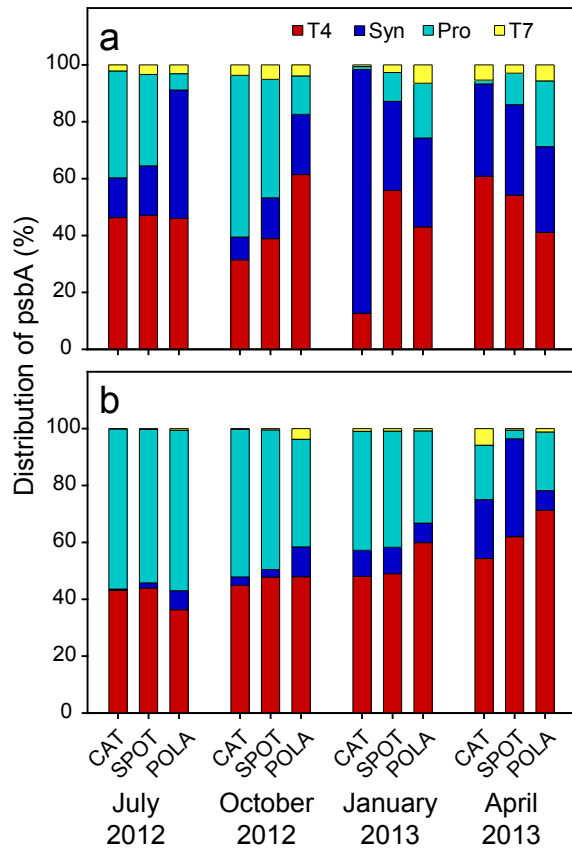147    closely with a different phage isolated on *Synechococcus* (sup. fig. 5).
148
149
150

**Figure 3**: Presence and activity of the assembled cyanophage and its potential hosts in October 2012: (**A**) Mean coverage (quartiles 1-3) of assembled cyanophage (**B**) OTU relative abundance by 16S-rRNA of the two most abundant cyanobacteria OTUs in order: *Prochlorococcus* in DNA, *Prochlorococcus* in RNA, *Synechococcus* in DNA, *Synechococcus* in RNA. Note the near-absence of *Prochlorococcus* in POLA, in contrast to *Synechococcus* and the phage, leading us to infer the phage infects *Synechococcus*.

Because viruses and hosts both code for photosynthetic functions, a comparison of viral and host-coded contributions to activity is possible. Sharon et al.[22] previously showed viral *psbA* gene can outnumber cyanobacterial *psbA* genes in metagenomes from the Mediterranean, and showed viral gene expression is evident. We extended this to quantitatively partition gene expression into bacterial contribution from *Synechococcus* and *Prochlorococcus* and viral contribution from cyanomyoviruses and cyanopodoviruses, as evident from HMM-placed translated reads onto our PS-II D1 phylogenetic tree. We found *psbA* transcripts of T4-like cyanomyovirus origin generally accounted for roughly 50% of cyanobacterial and cyanophage *psbA* transcripts. *Prochlorococcus* transcripts were almost always comparable to the T4-like contribution. On several occasions, the viral version exceeded the cyanobacterial version in read count (fig. 4).

We can roughly estimate the proportion of infected cyanobacteria from our *psbA* data and compare it to previously published estimates. For cyanobacteria in marine systems, the highest estimates of infection are roughly 50-60% infected at any given time[2,17,23,24]. One consideration when calculating the proportion of infected cyanobacteria is that during host infection, the number of phage mRNA of *psbA* increases quickly during early infection until it becomes the exclusive source of *psbA* transcripts in the cell[20,21]. Another consideration is that, regardless their source, host or virus, the abundance of *psbA* transcripts is comparable in infected and uninfected

178    cells[23]. What we observe in the sample is a comparable contribution of T4-like phages and
179    cyanobacteria (fig. 5 D) at a ratio of 1.2±0.6 (mean ± standard deviation) phage/cyanobacteria,
180    which suggests that on average about half of the cyanobacteria are infected. This is in accordance
181    with the high end of published estimates, confirming that infection is an important part of
182    cyanobacterial ecology.
183



184
185
186    **Figure 4:** Distribution of psbA of T4-like phages, *Synechococcus*, *Prochlororoccus*, and T7-like
187    phages in (A) metagenomes and (B) metatranscriptomes.

188
189    In both metagenomes and metatranscriptomes, there is minor consistent recruitment to T7-like
190    cyanopodovirus psbA. However, in every sample the contribution of T7-like cyanopodoviruses
191    was very low compared to that of T4-like cyanomyoviruses. This could be due to the more
192    specific host range reported for cyanopodoviruses compared to cyanomyoviruses[25-27]. As T4-like
193    and T7-like cyanophages are reported to be strictly lytic[28], their presence in metagenomes results
194    from late infection genomic copies or virions within host cells, pseudolysogeny or phages that
195    adsorbed to cells or particles.
196    Extending metatranscriptomics methods as recently applied to marine eukaryotic viral
197    infection[13,29,30], we show the power of multiple approaches to track viral infection and dynamics
198    within the broad picoplankton community, using metatranscriptomes of the cellular fraction,
199    with particular examples in the cyanobacteria. Use of marker genes is especially important to
200    study viruses with many close relatives in the same environment (whose contigs assemble
201    poorly), whereas assemblies are useful for tracking ephemeral, more clonal viruses. The

202  observed infection dynamics can sometimes be used in combination with microbial community
203  structure and viral marker genes found within contigs to deduce a host. Use of metagenomes and
204  metatranscriptomes provides an insight into quantifiable viral contribution to photosynthesis and
205  to estimating the fraction of infected cyanobacteria.
206
207

208  **Methods**
209  *Sample collection*
210  Surface seawater was collected by bucket on 7/15/2012, 10/19/2012, 1/9/2013 and 4/24/2013 in
211  three locations: The Port of Los Angeles (33°42.75'N 118°15.55'W), the San Pedro Ocean Time-
212  series (33°33.00'N 118°24.01'W) and Two Harbors, Santa Catalina Island (33°27.18'N
213  118°28.51'W). Duplicate samples of 20 liters were filtered in each location through an 80 μm
214  mesh followed by a glass fiber syringe prefilter (Gelman, 4523) which collected the >1 μm size
215  fraction and a 0.2 μm PES Sterivex filter (Millipore, SVGPB1010) which collected the free-
216  living size fraction. RNAlater (Thermo-Fisher, AM7020) was added to each filter and filters
217  were flash frozen no more than 5 minutes post-filtration.
218  *Library preparation*
219  DNA and RNA were extracted simultaneously from Sterivex filters by bead-beating followed by
220  an AllPrep kit (Qiagen, 80204). An internal standard (ERCC RNA Spike-In Mix, Thermo-Fisher
221  4456740) was added into the lysate after bead-beating for quality assurance. RNA was enriched
222  for mRNA with RiboZero (Illumina, MRZB12424). Resulting mRNA was reverse transcribed
223  using SuperScript-III (Invitrogen, 18080-051). DNA and cDNA were sheared with Covaris m2
224  and size-selected for products larger than 300 bp. RNA libraries were prepared and barcoded
225  using NEBNext Ultra Directional RNA library Prep Kit for Illumina (E74205). DNA libraries
226  were prepared and barcoded with Ovation UltraLow Library Prep V2 (Nugen, 0344).
227  Metagenomes were sequenced on Illumina HiSeq 2x125 bp or 2x150 bp. Metatranscriptomes
228  were sequenced on Illumina HiSeq 2x250 bp.
229  *Read processing and assembly*
230  Raw metagenomics and metatranscriptomics reads were quality trimmed and filtered with
231  Trimmomatic version 0.33 with parameters LEADING:20 TRAILING:20
232  SLIDINGWINDOW:15:25[31]. Metatranscriptomic reads were merged with PEAR[32], using the
233  default settings and residual ribosomal reads as well as the internal standard were removed
234  informatically. Merged reads from each sample separately were assembled with Megahit.
235  Contigs smaller than 2kbp from all samples were co-assembled with Newbler[33] version 2.9
236  (Roche) (minimum overlap 40bp minimum id 99%) and contigs larger than 2kbp from all
237  samples were co-assembled with minimus2[34] (minimum overlap 40bp minimum id 99%). Only
238  contigs larger than 5 Kbp were further analyzed.
239  *Identification and annotation of viral contigs*
240  Viral contigs were identified by VirSorter[35] using RefSeq on the CyVerse platform and only
241  contigs classified as category 1 or category 2 were considered. In addition, the contigs were
242  ranked using VirFinder[36] (rank >=0.95). Prodigal[37] was used to predict ORFs in those contigs,
243  and the amino acid sequences were searched against the nr database (August 12[th] 2016) using
244  blastp[38] and a maximum E-value $10^{-5}$. The annotations were used to verify viral contigs from the
245  VirFinder results. Contigs were verified to be non-chimeric by even recruitment.

246 Quality filtered metagenomic and metatranscriptomic reads were mapped back to these contigs
247 with Bowtie2 version 2.2.6 using the default settings and the expression patterns were identified
248 and visualized with Anvi'o[39] version 2.1.0.

*Microbial community composition analysis*

250 The V4-V5 regions of the 16S-rRNA coding gene were amplified from DNA and cDNA from all
251 samples using the 515-N-F and 926-R primers, and sequenced on an Illumina MiSeq 2x300 bp
252 (UC Davis genome center) along with a mock community as described in Parada et al.[40].
253 The ends of resulting reads were trimmed with PRINSEQ[41] to a quality score higher than 20. The
254 trimmed reads were merged with USEARCH7[42] allowing for 3 mismatches in the overlap region.
255 Retained assembled reads were clustered with mothur[43] version 1.38.0 according to the MiSeq
256 and classified with SILVA version 119. Bray-Curtis dissimilarity and dendrograms were
257 calculated and plotted with R package vegan[44].

*Analysis of PS-II D1 protein sequences*

259 A curated set of PS-II D1 amino acid sequences of myoviruses, podoviruses, cyanobacteria and
260 eukaryotes (chloroplast) from Pfam[45] and RefSeq release 80 was downloaded. All sequences of
261 marine viral PS-II D1 were retained in addition to sequences of bacterial and eukaryotic taxa that
262 were identified in the 16S-rRNA community composition. One of the assembled contigs
263 contained a psbA gene coding for PS-II D1. The translated amino acid sequences were added to
264 the set of proteins.
265 Merged reads from the metatranscriptomes and unmerged forward reads from the metagenomes
266 were aligned with blastx[38] against this set demanding an e-value of $10^{-5}$. The reads that passed
267 the filter were translated using bioPython[46] into amino acids according to the reading frame
268 indicated by the blastx start and end values.
269 Following the protocol used in Ignacio-Espinoza et al.[47] total of 158 sequences were aligned with
270 mafft[48] version 7.305b with parameters set to globalpair, gap open penalty 1.5, gap extension
271 penalty 0.5 and scoring matrix BLOSUM30. Informative blocks were identified using Gblocks[49]
272 version 0.91b with a minimum block length 5, blocks represent at least half of the sequences and
273 allowing gaps (b3=50, b4=5, b5=h). The blocks were used to build a maximum likelihood
274 phylogenetic tree using RAxML[50] (best of 20 trees, gamma model and WAG substitution
275 matrix). A hidden Markov Model (HMM) of the same set was also built with hmmer 3.0[51]. The
276 translated metagenomics and metatranscriptomics amino acid sequences were searched using the
277 HMM and a threshold of e-value $10^{-5}$. A total of 190,928 translated metatranscriptomics reads
278 and 72,292 metagenomics reads from all samples remained after this step. Those reads were
279 locally aligned to the HMM using hmmer 3.0 function hmmalign and placed into the
280 phylogenetic tree using pplacer[52] version v1.1.alpha17 (sup. fig. 6).

*Analysis of gp23 protein sequences*

282 Metatranscriptomic and metagenomics reads were searched against a set of T4-like clusters of
283 orthologous groups (COGs) with an E-value threshold of $10^{-5}$. 89,768 metatranscriptomic reads
284 and 134,995 metagenomic reads were annotated as gp23. An HMM of gp23 was built as
285 described previously and translated reads were searched and placed with pplacer. The tree was
286 visualized by the Interactive Tree Of Life (iTOL)[53].

*Recruitment to phage genomes*

288 The four currently available full pelagiphage genomes were downloaded from NCBI and
289 concatenated with assembled viral contigs from metatranscriptomes the metagenomes as well as
290 with published cyanophage genomes downloaded from NCBI RefSeq. Metagenomic and
291 metatranscriptomics reads were searched against the genomes dataset with blastn default

292    settings. For metagenomes only hits longer than 100bp were retained, and for
293    metatranscriptomes only hits longer than 200bp. Hits were then plotted against the genomes
294    using R[54].

295    *Data availability*
296    All data can be found on EMBL-ENA under project number PRJEB12234. Raw
297    metatranscriptomics sequences accession numbers are ERS1864892-ERS1864903, and negative
298    control library sequences accession number is ERR2089009. Raw metagenomic sequences
299    accession numbers are ERS1869885-ERS1869896 and negative control accession number is
300    ERS1872073. Assembled viral contigs accession numbers are ERZ474118-ERZ474183.

301

309

310    **References**:
311    1.  Proctor, L. M., & Fuhrman, J. A. (1990). Viral mortality of marine bacteria and
312        cyanobacteria. *Nature*, *343*(6253), 60.
313    2.  Suttle, C. A. (2007). Marine viruses--major players in the global ecosystem. *Nature*
314        *reviews. Microbiology*, *5*(10), 801.
315    3.  Lima-Mendez, G., Faust, K., Henry, N., Decelle, J., Colin, S., Carcillo, F., ... & Bittner,
316        L. (2015). Determinants of community structure in the global plankton
317        interactome. *Science*, *348*(6237), 1262073.
318    4.  Brum, J. R., Ignacio-Espinoza, J. C., Roux, S., Doulcier, G., Acinas, S. G., Alberti, A., ...
319        & Gorsky, G. (2015). Patterns and ecological drivers of ocean viral
320        communities. *Science*, *348*(6237), 1261498.
321    5.  Paez-Espino, D., Eloe-Fadrosh, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann,
322        M., Mikhailova, N., ... & Kyrpides, N. C. (2016). Uncovering Earth's
323        virome. *Nature*, *536*(7617).
324    6.  Nishimura, Y., Watai, H., Honda, T., Mihara, T., Omae, K., Roux, S., ... & Sullivan, M.
325        B. (2017). Environmental viral genomes shed new light on virus-host interactions in the
326        ocean. *mSphere*, *2*(2), e00359-16.
327    7.  Chow, C. E. T., & Fuhrman, J. A. (2012). Seasonality and monthly dynamics of marine
328        myovirus communities. *Environmental microbiology*, *14*(8), 2171-2183.
329    8.  Adriaenssens, E. M., & Cowan, D. A. (2014). Using signature genes as tools to assess
330        environmental viral ecology and diversity. *Applied and environmental*
331        *microbiology*, *80*(15), 4470-4480.
332    9.  Miranda, J. A., Culley, A. I., Schvarcz, C. R., & Steward, G. F. (2016). RNA viruses as
333        major contributors to Antarctic virioplankton. *Environmental microbiology*, *18*(11),
334        3714-3727.
335    10. Ottesen, E. A., Young, C. R., Eppley, J. M., Ryan, J. P., Chavez, F. P., Scholin, C. A., &
336        DeLong, E. F. (2013). Pattern and synchrony of gene expression among sympatric marine

337    microbial populations. *Proceedings of the National Academy of Sciences*, *110*(6), E488-
338    E497.
339  11. Ottesen, E. A., Young, C. R., Gifford, S. M., Eppley, J. M., Marin, R., Schuster, S. C., ...
340    & DeLong, E. F. (2014). Multispecies diel transcriptional oscillations in open ocean
341    heterotrophic bacterial assemblages. *Science*, *345*(6193), 207-212.
342  12. Jia, Y., Shan, J., Millard, A., Clokie, M. R., & Mann, N. H. (2010). Light-dependent
343    adsorption of photosynthetic cyanophages to Synechococcus sp. WH7803. *FEMS*
344    *microbiology letters*, *310*(2), 120-126.
345  13. Moniruzzaman, M., Wurch, L. L., Alexander, H., Dyhrman, S. T., Gobler, C. J., &
346    Wilhelm, S. W. (2017). Virus-host relationships of marine single-celled eukaryotes
347    resolved from metatranscriptomics. *Nature Communications*, *8*.
348  14. Needham, D. M., Sachdeva, R., & Fuhrman, J. A. (2017). Ecological dynamics and co-
349    occurrence among marine phytoplankton, bacteria and myoviruses shows microdiversity
350    matters. *The ISME Journal*.
351  15. Awad, S., Irber, L., & Brown, C. T. (2017). Evaluating Metagenome Assembly on a
352    Simple Defined Community with Many Strain Variants. *bioRxiv*, 155358.
353  16. Martinez-Hernandez, F., Fornas, O., Gomez, M. L., Bolduc, B., de la Cruz Pena, M. J.,
354    Martínez, J. M., ... & Sullivan, M. B. (2017). Single-virus genomics reveals hidden
355    cosmopolitan and abundant viruses. *Nature Communications*, *8*.
356  17. Deng, L., Ignacio-Espinoza, J. C., Gregory, A. C., Poulos, B. T., Weitz, J. S.,
357    Hugenholtz, P., & Sullivan, M. B. (2014). Viral tagging reveals discrete populations in
358    Synechococcus viral genome sequence space. *Nature*, *513*(7517), 242.
359  18. Williamson, S. J., Allen, L. Z., Lorenzi, H. A., Fadrosh, D. W., Brami, D., Thiagarajan,
360    M., ... & Venter, J. C. (2012). Metagenomic exploration of viruses throughout the Indian
361    Ocean. *PLoS One*, *7*(10), e42047.
362  19. Hurwitz, B. L., & U'Ren, J. M. (2016). Viral metabolic reprogramming in marine
363    ecosystems. *Current opinion in microbiology*, *31*, 161-168.
364  20. Lindell, D., Jaffe, J. D., Johnson, Z. I., Church, G. M., & Chisholm, S. W. (2005).
365    Photosynthesis genes in marine viruses yield proteins during host
366    infection. *Nature*, *438*(7064), 86.
367  21. Clokie, M. R., Shan, J., Bailey, S., Jia, Y., Krisch, H. M., West, S., & Mann, N. H.
368    (2006). Transcription of a 'photosynthetic'T4-type phage during infection of a marine
369    cyanobacterium. *Environmental Microbiology*, *8*(5), 827-835.
370  22. Sharon, I., Tzahor, S., Williamson, S., Shmoish, M., Man-Aharonovich, D., Rusch, D. B.,
371    ... & Adir, N. (2007). Viral photosynthetic reaction center genes and transcripts in the
372    marine environment. *The ISME journal*, *1*(6), 492.
373  23. Proctor, L. M., & Fuhrman, J. A. (1990). Viral mortality of marine bacteria and
374    cyanobacteria. *Nature*, *343*(6253), 60.
375  24. Wommack, K. E., & Colwell, R. R. (2000). Virioplankton: viruses in aquatic
376    ecosystems. *Microbiology and molecular biology reviews*, *64*(1), 69-114.
377  25. Sullivan, M. B., Waterbury, J. B., & Chisholm, S. W. (2003). Cyanophages infecting the
378    oceanic cyanobacterium Prochlorococcus. *Nature*, *424*(6952), 1047.
379  26. Millard, A. D., & Mann, N. H. (2006). A temporal and spatial investigation of
380    cyanophage abundance in the Gulf of Aqaba, Red Sea. *Journal of the Marine Biological*
381    *Association of the United Kingdom*, *86*(3), 507-515.

382  27. Wang, K., & Chen, F. (2008). Prevalence of highly host-specific cyanophages in the
383       estuarine environment. *Environmental microbiology*, *10*(2), 300-312.
384  28. Martin, E., & Benson, R. (1988). Phages of cyanobacteria. *The bacteriophages*, *2*, 607-
385       645.
386  29. Dupont, C. L., McCrow, J. P., Valas, R., Moustafa, A., Walworth, N., Goodenough, U.,
387       ... & Mann, E. (2015). Genomes and gene expression across light and productivity
388       gradients in eastern subtropical Pacific microbial communities. *The ISME journal*, *9*(5),
389       1076.
390  30. Allen, L. Z., McCrow, J. P., Ininbergs, K., Dupont, C. L., Badger, J. H., Hoffman, J. M.,
391       ... & Venter, J. C. (2017). The Baltic Sea Virome: Diversity and Transcriptional Activity
392       of DNA and RNA Viruses. *mSystems*, *2*(1), e00125-16.
393  31. Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for
394       Illumina sequence data. *Bioinformatics*, *30*(15), 2114-2120.
395  32. Zhang, J., Kobert, K., Flouri, T., & Stamatakis, A. (2013). PEAR: a fast and accurate
396       Illumina Paired-End reAd mergeR. *Bioinformatics*, *30*(5), 614-620.
397  33. Margulies, M., Egholm, M., Altman, W. E., Attiya, S., Bader, J. S., Bemben, L. A., ... &
398       Dewell, S. B. (2005). Genome sequencing in open microfabricated high density picoliter
399       reactors. *Nature*, *437*(7057), 376.
400  34. Sommer, D. D., Delcher, A. L., Salzberg, S. L., & Pop, M. (2007). Minimus: a fast,
401       lightweight genome assembler. *BMC bioinformatics*, *8*(1), 64.
402  35. Roux, S., Enault, F., Hurwitz, B. L., & Sullivan, M. B. (2015). VirSorter: mining viral
403       signal from microbial genomic data. *PeerJ*, *3*, e985.
404  36. Ren, J., Ahlgren, N. A., Lu, Y. Y., Fuhrman, J. A., & Sun, F. (2017). VirFinder: a novel
405       k-mer based tool for identifying viral sequences from assembled metagenomic
406       data. *Microbiome*, *5*(1), 69.
407  37. Hyatt, D., Chen, G. L., LoCascio, P. F., Land, M. L., Larimer, F. W., & Hauser, L. J.
408       (2010). Prodigal: prokaryotic gene recognition and translation initiation site
409       identification. *BMC bioinformatics*, *11*(1), 119.
410  38. Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., &
411       Madden, T. L. (2009). BLAST+: architecture and applications. *BMC
412       bioinformatics*, *10*(1), 421.
413  39. Eren, A. M., Esen, Ö. C., Quince, C., Vineis, J. H., Morrison, H. G., Sogin, M. L., &
414       Delmont, T. O. (2015). Anvi'o: an advanced analysis and visualization platform for
415       'omics data. *PeerJ*, *3*, e1319.
416  40. Parada, A. E., Needham, D. M., & Fuhrman, J. A. (2016). Every base matters: assessing
417       small subunit rRNA primers for marine microbiomes with mock communities, time series
418       and global field samples. *Environmental microbiology*, *18*(5), 1403-1414.
419  41. Schmieder, R., & Edwards, R. (2011). Quality control and preprocessing of metagenomic
420       datasets. *Bioinformatics*, *27*(6), 863-864.
421  42. Edgar, R. C. (2010). Search and clustering orders of magnitude faster than
422       BLAST. *Bioinformatics*, *26*(19), 2460-2461.
423  43. Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., ...
424       & Sahl, J. W. (2009). Introducing mothur: open-source, platform-independent,
425       community-supported software for describing and comparing microbial
426       communities. *Applied and environmental microbiology*, *75*(23), 7537-7541.

427  44. Oksanen, J., Kindt, R., Legendre, P., O'Hara, B., Stevens, M. H. H., Oksanen, M. J., &
428      Suggests, M. A. S. S. (2007). The vegan package. *Community ecology package*, *10*, 631-
429      637. http://vegan.r-forge.r-project.org
430  45. Finn, R. D., Coggill, P., Eberhardt, R. Y., Eddy, S. R., Mistry, J., Mitchell, A. L., ... &
431      Salazar, G. A. (2016). The Pfam protein families database: towards a more sustainable
432      future. *Nucleic acids research*, *44*(D1), D279-D285.
433  46. Cock, P. J., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., ... & De
434      Hoon, M. J. (2009). Biopython: freely available Python tools for computational
435      molecular biology and bioinformatics. *Bioinformatics*, *25*(11), 1422-1423.
436  47. Ignacio-Espinoza, J. C., & Sullivan, M. B. (2012). Phylogenomics of T4 cyanophages:
437      lateral gene transfer in the 'core'and origins of host genes. *Environmental*
438      *microbiology*, *14*(8), 2113-2126.
439  48. Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software
440      version 7: improvements in performance and usability. *Molecular biology and*
441      *evolution*, *30*(4), 772-780.
442  49. Castresana, J. (2000). Selection of conserved blocks from multiple alignments for their
443      use in phylogenetic analysis. *Molecular biology and evolution*, *17*(4), 540-552.
444  50. Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-
445      analysis of large phylogenies. *Bioinformatics*, *30*(9), 1312-1313.
446  51. Johnson, L. S., Eddy, S. R., & Portugaly, E. (2010). Hidden Markov model speed
447      heuristic and iterative HMM search procedure. *BMC bioinformatics*, *11*(1), 431.
448  52. Matsen, F. A., Kodner, R. B., & Armbrust, E. V. (2010). pplacer: linear time maximum-
449      likelihood and Bayesian phylogenetic placement of sequences onto a fixed reference
450      tree. *BMC bioinformatics*, *11*(1), 538.
451  53. Letunic, I., & Bork, P. (2016). Interactive tree of life (iTOL) v3: an online tool for the
452      display and annotation of phylogenetic and other trees. *Nucleic acids research*, *44*(W1),
453      W242-W245.
454  54. R Core Team (2016) https://www.R-project.org/