

1 **Full Article**

2 **CRISPR spacers indicate preferential matching of specific**
3 **virio plankton genes**

4 Daniel J. Nasko^{1†}, Barbra D. Ferrell¹, Ryan M. Moore¹, Jaysheel D. Bhavsar¹, Shawn
5 W. Polson¹, K. Eric Wommack^{1*}
6

7 *1. Delaware Biotechnology Institute, University of Delaware, Newark, Delaware, USA*
8

9 *To whom correspondence should be addressed.

10 Corresponding Authors' Information

11 **Address: Delaware Biotechnology Inst., 15 Innovation Way, Newark, Delaware**
12 **19711**

13 **(Tel): (302) 831-4362**

14 **(Fax): (302) 831-3447**

15 **(E-mail): wommack@dbi.udel.edu**
16

17 † Current affiliation: Institute for Advanced Computer Studies, University of Maryland,
18 College Park, Maryland, USA
19
20
21
22
23
24

25 **Abstract**

26 Viral infection exerts selection pressure on marine microbes as viral-induced cell lysis
27 causes 20 to 50% of cell mortality resulting in fluxes of biomass into oceanic dissolved
28 organic matter. Archaeal and bacterial populations can defend against viral infection
29 using the CRISPR-Cas system which relies on specific matching between a spacer
30 sequence and a viral gene. If a CRISPR spacer match to any gene within a viral
31 genome is equally effective in preventing lysis, then no viral genes should be
32 preferentially matched by CRISPR spacers. However, if there are differences in
33 effectiveness then certain viral genes may demonstrate a greater frequency of CRISPR
34 spacer matches. Indeed, homology search analyses of bacterioplankton CRISPR
35 spacer sequences against viroplankton sequences revealed preferential matching of
36 replication proteins, nucleic acid binding proteins, and viral structural proteins. Positive
37 selection pressure for effective viral defense is one parsimonious explanation for these
38 observations. CRISPR spacers from viroplankton metagenomes preferentially matched
39 methyltransferase and phage integrase genes within viroplankton sequences. These
40 viroplankton CRISPR spacers may assist infected host cells in defending against
41 competing phage. Analyses also revealed that half of the spacer-matched viral genes
42 were unknown and that some genes matched several spacers and some spacers
43 matched multiple genes, a many-to-many relationship. Thus, CRISPR spacer matching
44 may be an evolutionary algorithm, agnostically identifying those genes under stringent
45 selection pressure for sustaining viral infection and lysis. Investigating this subset of
46 viral genes could reveal those genetic mechanisms essential to viral-host interactions
47 and provide new technologies for optimizing CRISPR defense in beneficial microbes.

48 **MAIN TEXT**

49 Between 20 and 50% of microbial mortality within marine systems results from viral
50 infection and lysis. As a consequence, these processes are critical in driving carbon
51 and nutrient cycles within the sea (1, 2). In response to the substantial pressure of viral
52 predation, a number of sophisticated defense systems have evolved within cellular
53 microbial hosts including: alteration of cell surface receptors, production of extracellular
54 polysaccharides (3), restriction modification systems (4), and the clustered regularly
55 interspaced short palindromic repeat (CRISPR) system. Of these systems, the CRISPR
56 system is perhaps the most adaptable and specific, acting as an acquired immune
57 system in Bacteria and Archaea against bacteriophage and archaeal viruses
58 respectively, as well as other invading foreign DNA, such as plasmids (5). The
59 adaptability of the CRISPR system for targeting specific DNA regions for nuclease
60 digestion has been leveraged into a new and powerful approach for selective genome
61 editing within complex plant and animal genomes (6).

62 The CRISPR locus is comprised of CRISPR-associated (*cas*) genes and one or
63 more CRISPR sequence arrays consisting of a repeating pattern of different spacer
64 sequences and the same hairpin repeat sequence. It is the spacers that enable the
65 adaptable and gene-specific inactivating mechanism of the CRISPR system. Spacers
66 are short segments (26-72 base pairs (7)) of sequence that are homologous to phage or
67 plasmid DNA. Each spacer is flanked by comparably sized repeat sequences. The
68 repeats form a hairpin secondary structure and are conserved among bacterial and
69 archaeal species. The number of spacers in a CRISPR array varies from 2 to over 200

70 (7) and, interestingly, the position of a spacer in the array can provide an historical
71 timeline of viral host encounters (5).

72 After transcription, Cas proteins cleave repeats from the array transcript creating
73 small interfering CRISPR RNAs (crRNAs). The crRNAs are comprised of one spacer
74 flanked on either side by half a repeat. If a spacer sequence within a crRNA matches a
75 segment of an invading virus' genome, then the small interfering crRNA will target the
76 genomic DNA or RNA for destruction by the Cas proteins thus preventing viral
77 replication and ultimately cell mortality (8). Assuming that every gene a virus carries in
78 its genome is essential for successful infection and lysis, then, successful CRISPR
79 inactivation of any viral gene should prevent cell mortality from viral lysis. Given this
80 understanding of CRISPR defense against viral infection, we should expect no
81 preferential matches of viral genes to CRISPR spacer sequences. However, if there are
82 differences in the effectiveness of inactivating certain viral genes over others, then
83 certain viral genes may demonstrate a greater propensity to be matched by CRISPR
84 spacers. This hypothesis was addressed by identifying spacers within microbial and
85 viral metagenome sequence libraries and investigating whether subsets of viral genes
86 were preferentially matched by these CRISPR spacers.

87 CRISPR spacers offer a powerful tool for investigating phage-host interactions as
88 spacer sequences can link phage and host populations within complex microbial
89 communities (9, 10). For example, this approach was used to identify the microbial
90 hosts of unknown viral populations within the extreme environments of deep-sea
91 hydrothermal vents (11, 12). The biochemical mechanism controlling the selection of
92 protospacer sequences (i.e. candidate spacers from invading viral and plasmid DNA)

93 relies on a short DNA motif (usually 2-6 base pairs) directly adjacent to protospacer
94 sequences (protospacer adjacent motif or PAM) (13)-(14). Because the PAM is a short
95 sequence, these motifs can be common within a viral genome and thus, the PAM alone
96 does not necessarily predispose particular viral genes as possible protospacer targets.
97 However, positive selection for more effective viral resistance would mean that certain
98 subsets of viral genes are preferentially represented as targets of CRISPR spacers
99 within natural viroplankton communities. Information on viral genes preferentially
100 matched by CRISPR spacers could indicate those viral genes most critical to successful
101 viral replication and lysis. Given that the function of most viral genes is unknown (15),
102 information on preferential spacer targeting could provide clues as to the subset of
103 unknown viral genes that are under stringent selection for successful infection and host
104 cell lysis. Fundamental information on the CRISPR susceptibility of particular viral
105 genes could be leveraged to engineer more effective phage resistance in beneficial
106 microbes.

107 Spacers can be identified within DNA sequence libraries based on their
108 characteristic repeat-spacer pattern within a CRISPR array. Several tools are currently
109 available for identifying CRISPR spacer arrays, however, these tools tend to have a
110 high false discovery rate of spacer sequences as repeat sequence arrays resembling
111 CRISPR spacer arrays are common within microbial genomes (16). To address this
112 shortcoming CASC (CASC Ain't Simply CRT) was developed as a discovery tool
113 capable of validating the accuracy of CRISPR spacer predictions. CASC employs a
114 modified version of the CRISPR Recognition Tool (CRT) (16) to identify putative
115 CRISPR arrays followed by novel heuristics (search for known repeats, spacer size

116 distribution check) to examine and validate each putative CRISPR array. CASC is able
117 to run in an exploratory (liberal) mode, as well as a stricter (conservative) mode in which
118 identified arrays must contain known repeat sequences or Cas protein genes near the
119 array.

120 After validation, CASC was used to identify CRISPR spacers within large
121 collections of marine microbial metagenome sequence data from the Global Ocean
122 Sampling (GOS) and *Tara* Oceans expeditions (17, 18). These spacers were then used
123 to examine phage-host interactions throughout the global ocean and identify common
124 genetic vulnerabilities among viral populations exploited by marine prokaryotes to
125 defend against viral infection.

126

127 **RESULTS**

128 **CASC validation with artificial data**

129 Two artificial metagenomes were created to simulate Illumina reads and
130 pyrosequencing reads. Both of these metagenomes were comprised of the same ten
131 bacterial genomes: five genomes containing CRISPR arrays and five genomes without
132 CRISPR arrays (see methods section).

133 The simulated Illumina sequence reads (150 bp, paired-end) were assembled
134 with SPAdes (19) and produced ca. 1,800 contigs (mean length of 17,700 bp). Only
135 one of the ten genomes (*C. trancomatis* F/SW5) was completely assembled into one
136 contiguous sequence. Although the remaining genomes were fragmented into many
137 contigs, the known CRISPR arrays were represented in the assembled dataset. The

138 second artificial metagenome was composed of ca. 1 million pyrosequencing reads
139 (450 bp) that were directly analyzed without assembly.

140 Each CRISPR algorithm evaluated (CASC, CRT, PILER-CR (20), and CRISPR
141 Finder (21)) performed better in terms of sensitivity (ability to detect spacer loci) and
142 precision (ability to detect only valid spacer loci) when searching for spacers within
143 assembled contigs from Illumina sequence libraries as opposed to pyrosequencing
144 reads (Tables S2 and S3). CASC's validation steps, which remove potentially spurious
145 CRISPR predictions, resulted in more accurate CRISPR spacer predictions (Illumina
146 contigs precision = 1.0; pyrosequencing reads precision = 0.82) than all of the other
147 tools that were evaluated.

148

149 **Spacer predictions in GOS metagenomes**

150 The GOS reads dataset provided spacers from a broad geographic cross-section
151 of bacterioplankton communities. Because the GOS sequence reads averaged 915
152 nucleotides in length it was possible to search for CRISPR arrays within unassembled
153 reads. CASC (in liberal mode) was used to search for CRISPR spacers in all read
154 sequences from GOS. CASC identified 12,606 CRISPR spacers (>99% did not match
155 known spacers) contained in 2,686 arrays coming from 90% of all GOS sites (Additional
156 file 1). The site with the most spacers (13% of all spacers observed within the entire
157 GOS dataset) was GS033 (Punta Cormorant Lagoon, Floreana Island, Ecuador), which
158 was the most heavily sequenced site. The number of spacers found was normalized by
159 mega base pairs of reads sequenced at that site. Sites with the highest normalized
160 spacer abundance were often lakes or lagoons (seven of the top ten), with most having

161 more than two spacers per mega base pair of sequenced reads.

162 Nucleotide position histograms of the forward and reverse compliment direction of
163 each CRISPR repeat sequence were used as a means of post hoc testing of CRISPR
164 spacer arrays identified as “bona fide” and “non-bona fide” using CASC (liberal mode).
165 Repeats within bona fide CRISPR spacer arrays showed distinct positional nucleotide
166 signatures, whereas repeats within non-bona fide CRISPR array repeats showed no
167 discernible signature as each position had an equal occurrence of each nucleotide (Fig.
168 S2). The presence of a distinct positional nucleotide signature in the CASC bona fide
169 repeats was indicative of a collection of true and functioning repeat sequences within
170 the GOS data.

171

172 **Spacer predictions in *Tara Oceans* metagenomes**

173 *Tara Oceans* assembled contigs contained more than twice as many spacers
174 (29,879; 95% did not match known spacers) as the GOS reads (Additional file 2), likely
175 due to the greater sequencing depth and number of samples in the *Tara Oceans*
176 dataset. However, calculating the frequency of CRISPR spacers per mega base pair of
177 sequence data was confounded by the fact that these data were collected from
178 assembled contigs as opposed to single unassembled reads. To overcome this, read
179 recruitment information was obtained for each *Tara Oceans* contig which enabled
180 normalization of spacer abundance within the dataset (see methods). Between 15 and
181 71% of read bases were successfully recruited to contigs among the 178 *Tara Oceans*
182 microbial metagenomes (Additional File 2). The fraction of each library associated with
183 CRISPR spacers varied from 1×10^{-4} to 5×10^{-8} (Additional File 2).

184 After normalizing for sequencing effort, normalized spacer abundance (NSA)
185 within the *Tara* Oceans metagenomes showed a positive correlation with sample depth
186 (Pearson $r = 0.42$, p -value = $4e-9$) (Fig. 1). The sample with the highest normalized
187 spacer abundance was 122_MES_0.45-0.8, a mesopelagic sample having nearly 5,000
188 spacers per read Gbp recruited. Indeed, many of the samples with high NSA were from
189 the mesopelagic zone (21 of the top 30). NSA showed a positive correlation with GC
190 content as well (Pearson $r = 0.51$, p -value = $1e-13$), which was not surprising to see as
191 GC content also correlated strongly with depth (Pearson $r = 0.74$, p -value = $2e-16$).

192

193 **Linking CRISPR abundance to taxonomic composition of microbial communities**

194 Observed CRISPR spacer abundances in the global oceans were analyzed with
195 respect to the previously reported taxonomic composition of prokaryotic plankton
196 communities within *Tara* Oceans metagenomes (22). Nearly 75% of archaeal 16S
197 rDNA operational taxonomic units (OTUs) exhibited a positive correlation with NSA.
198 Thus, as NSA increased, the abundance of archaeal OTUs was more likely to increase
199 than decrease. In contrast only 50% of bacterial OTUs exhibited a positive correlation
200 with NSA, meaning that as NSA increased the abundance of bacterial OTUs was
201 equally likely to increase or decrease (p -value = $1.1e-14$). Additionally, there was a
202 positive correlation between NSA and Bray-Curtis dissimilarity, an index to assess
203 microbial community similarity (Mantel $r = 0.30$, p -value = 0.01). Thus, the greater the
204 compositional differences between prokaryotic plankton communities the greater the
205 difference in their NSA values.

206 At varying depth zones, the SAR clades within the Alphaproteobacteria sub-phyla
207 were consistently among the most negatively correlating OTUs with respect to NSA
208 (Fig. 2). Interestingly, some taxa with OTUs that negatively correlated with NSA also
209 had OTUs that positively correlated with NSA. In general, the percentage of OTUs with
210 significant positive correlations to NSA increased with depth (Surface = 2.2%, deep
211 chlorophyll maximum (DCM) = 6.6%, Mesopelagic = 7.5%), while the percentage of
212 OTUs with negative correlations to NSA remained fairly steady with the exception of the
213 DCM (Surface = 0.45%, DCM = 0.01%, Mesopelagic = 0.46%).

214

215 **Some viral genes are more likely to become spacers**

216 Matching a CRISPR spacer from a metagenome to a viral gene target (VGT) is
217 challenging because: (i) the collection of known reference viral genomes poorly
218 represents environmental viruses (especially aquatic viruses); (ii) viral genes mutate
219 rapidly; and (iii) the short length of spacer sequences means that even alignments with
220 a high percent identity match may have high BLASTn E values (Expect Values). To
221 address these challenges, a large database of virome sequences comprising 206
222 aquatic viral metagenomes and totaling ca. eight giga base pairs (Gbp) of sequence
223 data (65 *Tara* Oceans assembled viromes, 141 unassembled public viromes) was
224 collected. All microbial spacers found in the GOS and *Tara* Oceans datasets were
225 searched against the virome database with BLASTn (E value $\leq 1e-1$, word size 7) to
226 identify matches between spacers and candidate VGTs. Nucleotide open reading
227 frames (ORFs) were predicted only for virome sequences with a spacer match, allowing

228 for the detection of spacers that spanned two adjacent ORFs, which proved to be rare
229 (3% of spacers).

230 A many-to-many relationship between CRISPR spacers and their candidate
231 VGTs was observed – i.e. some spacers showed homology to multiple virome ORFs,
232 and some virome ORFs showed hits from multiple spacers (Fig. 3). While the majority
233 of spacers were homologous to only one virome ORF (nearly 1,500 spacers, 45%),
234 there were a few spacers with homology to over 400 virome ORFs. These
235 cosmopolitan spacers often targeted less complex regions of structural proteins such as
236 short glutamic acid repeats within a portal protein.

237 In total, nearly a quarter (24%) of the CASC-identified (run this time in
238 conservative mode ensuring these were bona fide spacers) bacterioplankton spacers
239 had a nucleotide BLAST alignment with a virome open reading frame. Nearly half of the
240 translated viral ORFs (43%) had a match to a Phage SEED peptide (23), the majority of
241 which had an informative annotation; i.e. were not simply labeled “Phage protein” (Fig.
242 4).

243 All virome ORFs in the virome database were annotated using homology
244 information to Phage SEED proteins, enabling quantification of the expected frequency
245 of VGT annotations. In turn, annotation data was used to establish an expected
246 frequency for each viral gene annotation within the collection global ocean viromes.
247 Each of the top fifteen annotations assigned to VGTs were assigned more often than
248 expected (Table 1, Additional file 5). There were two exceptions that were targeted less
249 frequently than expected, genes encoding phage tail fiber (a set of structural proteins
250 attached to the base of the tail, used in host recognition and attachment) and DNA

251 helicase (a motor protein that separates double-stranded nucleic acid). Overall, the
252 VGT ORFs had a higher rate of homology to Phage SEED peptides than would be
253 expected indicating that VGTs of CRISPR defense are among the better-known subset
254 of viral genes (expected 2,257 no-hits, observed 1,920).

255 The *Tara* Oceans microbial shotgun metagenomes and viromes provided a rich
256 set of spacer-to-virome ORF matches. However, instances of bacterioplankton spacers
257 matching ORFs within a virome collected from the same water sample were rare. More
258 frequently bacterioplankton spacers had matches to virome ORFs from viromes
259 collected several thousand miles away (Fig. 5). This was the case for bacterioplankton
260 metagenomes collected from surface and deep chlorophyll maximum water samples.

261

262 **Viruses encoding CRISPR spacer arrays**

263 Previous studies have shown that phages infecting marine bacteria can carry the
264 genetic elements of the CRISPR/Cas system (24, 25). Over 2,000 CRISPR spacers
265 were observed within the aquatic viromes. To determine if the virome spacers targeted
266 a different subset of viral genes than the bacterioplankton spacers, the virome spacers
267 were also assessed against the aquatic virome database, in the same way as the
268 bacterioplankton metagenome spacers.

269 A greater frequency of virome spacers had a match to virome ORFs than that
270 seen for bacterioplankton spacers (30% versus 24%). Additionally, more of these VGT
271 ORFs of virome spacers could be annotated with Phage SEED than the
272 bacterioplankton spacers (55% versus 43%) (Fig. 6, Additional file 6). Again, all of the
273 ORFs in the virome database were annotated with Phage SEED to establish an

274 expected frequency for each viral gene annotation in the global oceans. Among the
275 informative annotations (annotations that were not simply “Phage protein”)
276 methyltransferase was targeted 21 times more often than expected (expected ca. 5
277 annotations, observed 100) by viral spacers, whereas microbial spacers targeted
278 methyltransferase only 4 times more often than expected. Indeed, methyltransferase
279 was among several gene targets that are differentially targeted between microbial and
280 virome spacers, including integrase and antitermination protein Q.

281

282 **DISCUSSION**

283 By and large, the focus of work investigating CRISPR as a microbial defense
284 strategy has been to determine the biochemical mechanisms behind spacer acquisition
285 and maintenance within bacterial (26) and archaeal (27) taxa. As a consequence these
286 studies have been conducted in model organisms within experimental laboratory
287 systems (28, 29), with some exceptions (30). Here we investigated the diversity and
288 frequency of unknown CRISPR/Cas systems within the global ocean, an approach that
289 broadly accounted for the influence of environmental selective pressures on the
290 acquisition and maintenance of CRISPR spacers. These investigations revealed that
291 particular subsets of virioplankton genes are highly targeted by the CRISPR defense
292 system of bacterioplankton and that there is a many-to-many relationship of spacers to
293 virioplankton genes.

294 Deeply sequenced shotgun bacterioplankton metagenomes enabled the search
295 for novel CRISPR spacers across a wide geographic range of aquatic environments.
296 Increasing sequence read lengths and yields from next generation sequencers have

297 enabled modern assembly algorithms to better resolve the repeat-rich CRISPR locus
298 (31) as seen through the high yield of CRISPR spacers in the *Tara* Oceans dataset.
299 Testing indicated that the addition of quality control heuristics in CASC provided a more
300 reliable set of CRISPR spacers than other CRISPR-finding algorithms.

301 With the rich set of CRISPR spacers mined directly from the environment it is
302 possible to compare our findings to those obtained through mathematical theory and
303 single-organism model systems. Normalized spacer abundance positively correlated
304 with sample depth indicating that the CRISPR/Cas system is an important defense
305 strategy for deep-sea bacterial and archaeal populations. The concentrations of hosts
306 and viruses is known to decrease with depth in the ocean (32), thus, this observation
307 agrees with previous work demonstrating that inducible immunity (i.e. CRISPR) is
308 preferred in conditions where the concentrations of host and virus are low (33). Not
309 only was NSA generally lower at the surface, where concentrations of hosts and viruses
310 tend to be greater, there were also several surface water bacterioplankton taxa that
311 exhibited strong negative correlations with NSA; chief among them were taxa within the
312 abundant SAR11 clade (Pelagibacterales) (34). This may be further evidence of the
313 limited effectiveness of CRISPR/Cas defense in competitive environments, as SAR11
314 members (notorious defense specialists) appear to favor other mechanisms of
315 bacteriophage resistance (e.g. cryptic escape (35)) rather than CRISPR/Cas.

316 A protospacer is the 30-40 bp segment of a viral gene that is incorporated into a
317 CRISPR array as a spacer. A motif, adjacent to the protospacer, called the PAM
318 (protospacer adjacent motif) is essential to the spacer acquisition machinery (14) in
319 Type I and II CRISPR/Cas systems. However considering the short and often

320 degenerate nature of PAMs (e.g. 2 bp, 16-fold degenerate (36)), hundreds to thousands
321 of potential PAM sites can exist within a viral genome. Thus, while the PAM plays a role
322 in determining the site within a viral gene that becomes a protospacer it remains
323 uncertain what, if anything, contributes to the retention of certain spacers within the
324 array in a natural system. Given the commonality of PAMs within viral genomes, the
325 most parsimonious explanation for the observed selection of particular VGTs within
326 viroplankton metagenomes is positive selection pressure for effective viral defense.
327 The CRISPR spacers observed within the bacterioplankton metagenomes were
328 maintained because they were the most successful in minimizing the damaging impacts
329 of viral infection and lysis on bacterioplankton populations. These data also provide
330 interesting insights concerning those genes that are most critical to the processes of
331 viral infection and lysis of bacterioplankton hosts.

332 In particular, these data show that there are conserved regions of potentially
333 evolutionary constrained viral genes that are targeted more often than expected by
334 CRISPR spacers from bacterioplankton populations. Genes encoding phage terminase
335 (enzymes that initiate DNA packaging by cutting the DNA concatemer),
336 methyltransferase (a family of enzymes that catalyze the transfer of a methyl group to
337 DNA or RNA), recombinase (enzymes that catalyze exchanges of nucleic acid within a
338 genome), and ssDNA-binding proteins (proteins that bind single-stranded DNA to
339 prevent it from re-forming a double-stranded molecule) were among the most
340 overtargeted genes within the viroplankton (Fig. 4). An inference from these
341 observations is that these viral genes are under particularly stringent selection pressure
342 which prevents the easy acquisition of point mutations that would ordinarily allow a viral

343 gene target to evade spacer recognition, the critical first step in CRISPR defense. Thus,
344 our analysis has pointed to particular gene functions that may have a heightened
345 importance to successful replication of marine viral populations.

346 The observation of thousands of spacers within nearly 20% of the viromes
347 surveyed (38 of 206) indicated a high prevalence of CRISPR-carrying viruses. The
348 impact of CRISPR-carrying viral populations in natural microbial communities may be
349 greater than expected. The frequent observation of virome spacers supports the recent
350 finding that cyanophages have been shown to carry CRISPR arrays and perhaps
351 transfer the arrays between related cyanobacteria to offer infection resistance from
352 competing phage (25). An enrichment in viral spacers targeting methyltransferase and
353 integrase genes may indicate that viral CRISPR arrays aid the host in targeting
354 competing temperate phage.

355 Interestingly, CRISPR spacers from bacterioplankton metagenomes targeted
356 certain genes less frequently than expected such as phage tail fiber genes. The
357 relatively simple structure of phage tail fiber protein would indicate a less stringent
358 selective pressure at the coding level, implying a greater opportunity for tail fiber gene
359 diversity. Indeed, phage tail fiber genes have been shown to not only be hypervariable,
360 but also undergo targeted hypervariation by retroelements in order to expand viral host
361 range (38, 39). Additionally, viral ORFs targeted by CRISPR spacers were less likely to
362 have an unassigned function than expected (actual unassigned functions = 598,
363 expected = 699) indicating CRISPR-targeted viral genes are more likely to have a
364 known functional role as opposed to non-targeted genes (Fig. 4 and Table 1).
365 Nevertheless, nearly half (41%) of these CRISPR-targeted viral genes were unknown

366 and would be considered viral genetic “dark matter” (40). This subset of CRISPR-
367 targeted but unknown viral “dark matter” genes likely play an important role in infection
368 and lysis processes.

369 Spacers matched virome ORFs in a many-to-many relationship, indicating that
370 some spacers were capable of targeting several different virome ORFs and several
371 virome ORFs were targeted by multiple spacers. In the latter case, these viral genes
372 appear to be highly targeted by the CRISPR/Cas system (Fig. 3). Instances of virome
373 ORFs being targeted by multiple spacers suggests that these ORFs are under
374 especially stringent selection pressure and are thus less likely to evade CRISPR
375 interference through single nucleotide point mutations. The over-targeting of these
376 ORFs also indicates that they are critical to viral replication and are thus more effective
377 targets for bacterioplankton CRISPR immunity.

378 Interestingly, less than 1% of spacers from *Tara Oceans* microbial metagenomes
379 matched virome ORFs from the same site (Fig. 5). One potential explanation for this
380 observation is that spacers found in a given bacterioplankton metagenome have
381 successfully minimized the replication of targeted viral populations to a level below
382 detection within a virome library. This observation is consistent with previous studies of
383 Archaeal-dominated systems (41, 42) and emphasizes a potential challenge of using
384 CRISPR to link viruses with their hosts within a single environmental sample. The
385 analysis of paired microbial/viral metagenomes over time may provide interesting
386 perspectives, as it could be possible to observe spacers targeting viruses from past
387 samples.

388 This study analyzed a large collection of CRISPR spacers from microbial
389 populations throughout the global oceans and has provided evidence that particular viral
390 genes are preferentially targeted by the CRISPR/Cas system. The identification of
391 certain viral gene classes that are more likely to become CRISPR spacers indicates that
392 these genes represent a genetic vulnerability for viral populations and that these genes
393 are potentially under strict selective pressure for successful viral infection and lysis.
394 CRISPR spacers sequenced from the environment have shown to be useful in linking
395 microbial hosts to their viruses (43). Our findings also indicate that spacer sequences
396 can identify those viral genes that represent the points of greatest genetic vulnerability
397 for natural viral populations. In this way, CRISPR/Cas may be thought of as a living
398 “evolutionary algorithm” (a field of artificial intelligence, which mimics natural selection to
399 solve complex problems) to agnostically identify viral genes that are most vulnerable.
400 These genes may then be further explored for uses in biotechnology (e.g. preventing
401 phage infections in processes relying on bacterial fermentation) or analysis of phage
402 diversity (as they are likely conserved).

403

404 **METHODS**

405 **CASC Pipeline**

406 The CASC pipeline can be broadly divided into two parts (Fig. S1): (A)
407 preliminary search for putative CRISPR spacers and (B) validation of putative CRISPR
408 arrays by Cas protein homology, CRISPR repeat homology, and the statistical
409 characteristics of spacer sizes. The preliminary search for CRISPR arrays employs a
410 modified version of the CRT (16). Modifications included a reformatting of the search

411 output, improved handling of multi-FASTA files, and the ability to utilize multiple CPUs
412 to lessen computational run time. These modifications improved the ability of CRT to
413 analyze large metagenomic datasets. Putative CRISPR arrays are then validated and
414 deemed "bona fide" CRISPRs if any of the following conditions are met: (i) the
415 sequence containing the candidate CRISPR array has a BLASTx match (E value $\leq 1e$ -
416 12) to a known UniRef 100 Cas protein cluster (44), (ii) the candidate CRISPR repeat
417 had a BLASTn match (E value $\leq 1e$ -5, word size 4) to a known CRISPR repeat from the
418 CRISPRdb reference database (7), or (iii) the standard deviation of spacer length within
419 the candidate CRISPR array was less than or equal to two base pairs. CASC offers a
420 "conservative" and a "liberal" CRISPR validation mode. In conservative mode,
421 conditions (i) or (ii) must be met, while in liberal mode conditions (i), (ii), or (iii) may be
422 met. CASC is available on GitHub (<https://github.com/dnasko/CASC>).

423

424 **Simulated Metagenome Construction**

425 Two shotgun sequence simulations were generated using Grinder (ver. 0.5.0)
426 (45) for the purpose of validating CASC and assessing performance. Ten complete
427 bacterial genomes were selected for the simulated metagenomes (Table S1), five of
428 which contained CRISPR arrays. The first simulation generated 60 million paired-end
429 150 base pair Illumina reads (read_dist=150 normal 0; insert_dist=300;
430 mutation_dist=poly4) and the second simulation generated 1 million 454
431 pyrosequencing reads (read_dist=450 normal 50; mutation_dist=poly4).

432 The Illumina simulated read pairs were assembled using the St. Petersburg
433 genome assembler (SPAdes) version 3.5.0 using all default settings (19) with the

434 exception of bypassing the pre-assembly read error correction process. The 454
435 simulated reads were not assembled and CRISPRs were predicted directly from the
436 reads.

437

438 **Performance Validation**

439 The known CRISPR array positions in five of the ten genomes were used to
440 assess the performance (i.e. sensitivity and precision) of several CRISPR identification
441 algorithms. Alignment of the Illumina assembled contigs against the reference
442 genomes identified the position of each CRISPR locus on the contigs and indicated that
443 all spacers were successfully assembled. The alignment-generated CRISPR positions
444 on the contigs were then used as the known CRISPR array positions. CRISPR array
445 positions within the 454 reads were determined using the genome coordinates provided
446 by Grinder.

447 Several algorithms, including CASC version 2.5 and the default settings of
448 metaCRT (a version of the CRT modified by Rho and colleagues) (46), PILER-CR (ver.
449 1.06) (20), and CRISPR Finder (21), were used to predict CRISPR arrays from the
450 Illumina assembled contigs and 454 reads (Table S2 and Table S3). Predicted spacers
451 from each program were clustered with the set of known spacers using CD-HIT-EST
452 (ver. 4.6) (47). Those spacers clustering at 100% identity with a known spacer were
453 counted as a true positive.

454 To better measure the abundance of spacers in the simulated Illumina
455 metagenome a recruitment of the simulated Illumina reads to assembled SPAdes
456 contigs was performed using Bowtie2 (ver. 2.1.0) (48). Coverage of each spacer was

457 calculated using SAMtools (ver. 1.2-2-gf8a6274) (49) and used to estimate the number
458 of spacer copies present in the simulated Illumina metagenome.

459

460 **Spacer predictions in GOS and *Tara* Oceans microbial metagenomes**

461 The Global Ocean Sampling (GOS) and *Tara* Oceans expeditions sampled and
462 sequenced microbial DNA from across the world's oceans (17, 18). The GOS dataset
463 was ideally suited for CRISPR prediction as the long read technology used for
464 sequencing these libraries was capable of encoding intact CRISPR arrays (50), and this
465 dataset has been used in previous studies of CRISPR prediction from metagenomic
466 data (51, 52). GOS sequences were downloaded from iMicrobe (imicrobe.us) and
467 included the GOS I expedition, GOS Baltic Sea, and GOS Banyoles (Additional file 1).
468 CRISPR spacers were predicted from 157 GOS sequence libraries totaling ca. 39
469 million reads and containing ca. 21 Gbp of genomic DNA from microorganisms typically
470 between 0.1 and 0.8 μm in size (note that filter sizes ranged from 0.002 to 20 μm based
471 on sample site) with CRISPR calling in 'liberal' mode.

472 The *Tara* Oceans expedition was a global-scale oceanic study that sampled and
473 sequenced metagenomes from 67 sites (53). In addition to sampling nearly every site
474 at varying depths, several sites were processed with multiple filter sizes (ranging from
475 0.2 to 3.0 μm), including 54 sites with paired microbial and viral fractions, making the
476 *Tara* Oceans dataset ideal for linking bacterial spacers with their viral gene targets in
477 the viromes. *Tara* Oceans metagenomes were predominantly sequenced using Illumina
478 HiSeq (100 bp, paired-end reads). Because Illumina reads are too short for accurate
479 searches of spacer arrays, assembled contigs were used instead (ca. 58 million contigs

480 totaling 62 Gbp). *Tara* Oceans assembled contigs were obtained from the European
481 Nucleotide Archive (<http://www.ebi.ac.uk/ena/about/tara-oceans-assemblies>).

482 In addition to counting the number of spacers found within each *Tara* contig, it
483 was necessary to calculate the abundance of each spacer by recruitment of the original
484 library of unassembled Illumina reads to *Tara* contigs. The reads corresponding to each
485 assembly were downloaded from NCBI's Sequence Read Archive and recruited to their
486 assembled contigs using Bowtie2 (very sensitive local setting). Read coverage of each
487 spacer was calculated using SAMtools and used as a proxy for the number of copies of
488 each spacer.

489 To measure how novel these spacers were, the GOS and *Tara* Oceans spacers
490 were clustered with known spacers from the CRISPRDB at 98% identity using CD-HIT-
491 EST (7, 47).

492

493 **Microbial Community Profiles with Respect to CRISPR Abundance**

494 The *Tara* Oceans observed OTUs "16S OTU Table" from Sunagawa et al. (22)
495 was downloaded from <http://ocean-microbiome.embl.de/companion.html> and imported
496 into QIIME (54). OTUs occurring ≤ 2 times were filtered out and 100 jackknife
497 subsamples were created with 35,461 observations (90% of the smallest sample) in
498 each. The community similarity test was performed with `beta_diversity.py` using Bray-
499 Curtis. Per-OTU correlations were calculated for each depth zone after splitting the
500 BIOM file accordingly and using `observation_metadata_correlation.py`. Only correlations
501 with Pearson's $r \geq 0.3$ or ≤ -0.3 with $p\text{-value} \leq 0.05$ were considered significant.

502

503 **Identification of GOS and *Tara* Oceans Spacer Targets**

504 Putative CRISPR spacers from the GOS and *Tara* Oceans microbial
505 metagenomes were searched against *Tara* Oceans viromes (Additional file 3) and a
506 subset of publicly available aquatic viromes (Additional file 4) available on the Viral
507 Informatics Resource for Metagenome Exploration (VIROME, virome.dbi.udel.edu) (55)
508 to identify candidate viral gene targets. Only spacers found with CASC in conservative
509 mode were used for this analysis to reduce the likelihood of identifying spurious
510 spacers.

511 Sequence alignment cut-offs used in previous studies comparing microbial
512 spacers to virome genes have varied, both in stringency and cut-off metric, depending
513 on the aim of the study. When identifying host-phage interactions by linking specific
514 viral population(s) to CRISPR spacers/loci, more stringent cut-offs are applied, such as
515 requiring a 100% nucleotide identity alignment of ≥ 20 bp (11), or an alignment with no
516 more than one mismatch (56). Exploratory studies trying to link what, if any, similarities
517 exist between microbial spacers and virome genes have used more relaxed cut-offs,
518 such as E value $\leq 1e-3$ (10), or alignments containing up to 15 mismatches (57).

519 As the objective of this study was to determine if particular viral genes were more
520 likely to be targeted by the CRISPR system of marine bacterioplankton the latter, more
521 exploratory approach was used. Spacer sequences are highly diverse and hyper
522 variable, even between closely related species (58), making it challenging to identify
523 candidate viral gene targets at the nucleotide level. Thus, when searching for potential
524 viral gene targets in viromes some mismatches and gaps in the nucleotide alignment
525 were permitted using BLASTn (ver. 2.2.30+, E value $\leq 1e-1$, word size 7). This resulted

526 in 51% of high-scoring segment pairs (HSPs) with no mismatches and 89% of HSPs
527 with no gap openings (Fig. S3).

528 In this analysis some spacers matched CRISPR arrays within several viromes.
529 To limit these spurious matches, CASC (liberal mode) was used to identify putative
530 spacer arrays within the viromes. Subsequently, sequences containing an array were
531 removed from the aquatic virome database prior to the analysis to identify viral gene
532 targets.

533 Spacer sequences were searched against the virome database with BLASTn.
534 Virome sequences that aligned with spacers were then culled into a separate FASTA
535 file and open reading frames (ORFs) were predicted using MetaGene (59). ORFs were
536 predicted after the spacer search to detect any spacers that may have spanned virome
537 ORFs (a rare occurrence). Virome ORFs with a match to a spacer were translated and
538 searched against Phage SEED (version 01-May-2016) (<http://www.phantome.org>) using
539 BLASTp (ver. 2.2.30+, E value $\leq 1e-3$). Each ORF was annotated using the best
540 cumulative bit score, which is described in the next section.

541 Finally, great-circle distances between microbial metagenome spacers and VGTs
542 within viromes were calculated in R (60) using the geosphere package (61). Distance
543 distributions were rendered in violin plots using the R package vioplot.

544

545 **Annotating virome ORFs and calculating expectation**

546 Virome ORFs with a match to a spacer were translated and searched against
547 Phage SEED (version 01-May-2016) (<http://www.phantome.org>) using BLASTp (ver.
548 2.2.30+, E value $\leq 1e-3$). A virome ORF was annotated to be the gene function

549 producing the highest cumulative bit score. For example, if “ORF_1” hit ten Phage
550 SEED genes, eight of which were hits to phage protein and the total bit score of these
551 alignments was 50, while the two remaining hits were to terminases with a total bit score
552 of 100, then “ORF_1” would be assigned to terminase. ORF annotation counts were
553 generated for the virome ORFs matching microbial (Additional File 5) and virome
554 spacers (Additional File 6).

555 To put these counts in come context, all aquatic virome ORFs were run through
556 the same Phage SEED-based annotation pipeline. Counts for all virome ORFs were
557 tabulated and the frequency of occurrence for each gene type was calculated. The
558 expected number of genes to have matches to CRISPR spacers was calculated by
559 multiplying the total number of genes matching spacers by the frequency of that gene
560 being annotated in all aquatic viromes.

561

562 **Data Availability**

563 Scripts used in this analysis are available on GitHub (github.com/dnasko/CASC) under
564 the GNU General Purpose License.

565

566 Six datasets were used in this analysis. The first two were simulated metagenomic
567 datasets and are available at Zenodo (<http://doi.org/10.5281/zenodo.1650429>). The
568 second two datasets were shotgun metagenomic reads from the Global Ocean Survey
569 (GOS) and *Tara* Oceans survey. GOS sequences were downloaded from iMicrobe
570 (imicrobe.us) and included the GOS I expedition, GOS Baltic Sea, and GOS Banyoles
571 (Additional file 1). *Tara* Oceans assembled contigs were obtained from the European

572 Nucleotide Archive (<http://www.ebi.ac.uk/ena/about/tara-oceans-assemblies>). The fifth
573 dataset was a subset of publicly available aquatic viromes (Additional file 4) available on
574 the Viral Informatics Resource for Metagenome Exploration (VIROME,
575 virome.dbi.udel.edu). Finally, the *Tara Oceans* observed OTUs “16S OTU Table” from
576 Sunagawa et al. (22) was downloaded from [http://ocean-](http://ocean-microbiome.embl.de/companion.html)
577 [microbiome.embl.de/companion.html](http://ocean-microbiome.embl.de/companion.html).

578

579 REFERENCES

- 580 1. Weitz JS, Stock CA, Wilhelm SW, Bourouiba L, Coleman ML, Buchan A, Follows
581 MJ, Fuhrman JA, Jover LF, Lennon JT, Middelboe M, Sonderegger DL, Suttle CA,
582 Taylor BP, Frede Thingstad T, Wilson WH, Eric Wommack K. 2015. A multitrophic
583 model to quantify the effects of marine viruses on microbial food webs and
584 ecosystem processes. *ISME J* 9:1352–1364.
- 585 2. Poorvin L, Rinta-Kanto JM, Hutchins DA, Wilhelm SW. 2004. Viral release of iron
586 and its bioavailability to marine plankton. *Limnol Oceanogr* 49:1734–1741.
- 587 3. Labrie SJ, Samson JE, Moineau S. 2010. Bacteriophage resistance mechanisms.
588 *Nat Rev Microbiol* 8:317–327.
- 589 4. Dorman CJ. 2004. H-NS: a universal regulator for a dynamic genome. *Nat Rev*
590 *Microbiol* 2:391–400.
- 591 5. Sorek R, Kunin V, Hugenholtz P. 2008. CRISPR--a widespread system that
592 provides acquired resistance against phages in bacteria and archaea. *Nat Rev*
593 *Microbiol* 6:181–186.
- 594 6. Ran FA, Hsu PDP, Wright J, Agarwala V, Scott D a, Zhang F. 2013. Genome
595 engineering using the CRISPR-Cas9 system. *Nat Protoc* 8:2281–2308.
- 596 7. Grissa I, Vergnaud G, Pourcel C. 2007. The CRISPRdb database and tools to
597 display CRISPRs and to generate dictionaries of spacers and repeats. *BMC*
598 *Bioinformatics* 8:172.
- 599 8. Barrangou R, Fremaux C, Deveau H, Richards M, Moineau S, Romero DA,
600 Horvath P, Barrangou R, Fremaux C, Deveau H, Richards M. 2007. CRISPR
601 Provides Against Viruses Resistance Acquired in Prokaryotes. *Science* (80-)
602 315:1709–1712.
- 603 9. Andersson AF, Banfield JF. 2008. Virus population dynamics and acquired virus
604 resistance in natural microbial communities. *Science* 320:1047–50.
- 605 10. Berg Miller ME, Yeoman CJ, Chia N, Tringe SG, Angly FE, Edwards RA, Flint HJ,
606 Lamed R, Bayer EA, White BA. 2012. Phage-bacteria relationships and CRISPR
607 elements revealed by a metagenomic survey of the rumen microbiome. *Environ*
608 *Microbiol* 14:207–227.
- 609 11. Anderson RE, Brazelton WJ, Baross JA. 2011. Using CRISPRs as a

- 610 metagenomic tool to identify microbial hosts of a diffuse flow hydrothermal vent
611 viral assemblage. *FEMS Microbiol Ecol* 77:120–133.
- 612 12. Paez-Espino D, Chen IMA, Palaniappan K, Ratner A, Chu K, Szeto E, Pillay M,
613 Huang J, Markowitz VM, Nielsen T, Huntemann M, Reddy TBK, Pavlopoulos GA,
614 Sullivan MB, Campbell BJ, Chen F, McMahon K, Hallam SJ, Denev V, Cavicchioli
615 R, Caffrey SM, Streit WR, Webster J, Handley KM, Salekdeh GH, Tsesmetzis N,
616 Setubal JC, Pope PB, Liu WT, Rivers AR, Ivanova NN, Kyrpides NC. 2017.
617 IMG/VR: A database of cultured and uncultured DNA viruses and retroviruses.
618 *Nucleic Acids Res* 45:D457–D465.
- 619 13. Horvath P, Barrangou R. 2010. CRISPR/Cas, the immune system of bacteria and
620 archaea. *Science* 327:167–170.
- 621 14. Paez-Espino D, Morovic W, Sun CL, Thomas BC, Ueda K, Stahl B, Barrangou R,
622 Banfield JF. 2013. Strong bias in the bacterial CRISPR elements that confer
623 immunity to phage. *Nat Commun* 4:1430.
- 624 15. Rosario K, Breitbart M. 2011. Exploring the viral world through metagenomics.
625 *Curr Opin Virol* 1:289–297.
- 626 16. Bland C, Ramsey TL, Sabree F, Lowe M, Brown K, Kyrpides NC, Hugenholtz P.
627 2007. CRISPR recognition tool (CRT): a tool for automatic detection of clustered
628 regularly interspaced palindromic repeats. *BMC Bioinformatics* 8:209.
- 629 17. Yooseph S, Sutton G, Rusch DB, Halpern AL, Williamson SJ, Remington K, Eisen
630 JA, Heidelberg KB, Manning G, Li W, Jaroszewski L, Cieplak P, Miller CS, Li H,
631 Mashiyama ST, Joachimiak MP, van Belle C, Chandonia JM, Soergel DA, Zhai Y,
632 Natarajan K, Lee S, Raphael BJ, Bafna V, Friedman R, Brenner SE, Godzik A,
633 Eisenberg D, Dixon JE, Taylor SS, Strausberg RL, Frazier M, Venter JC. 2007.
634 The Sorcerer II Global Ocean Sampling expedition: expanding the universe of
635 protein families. *PLoS Biol* 5:e16.
- 636 18. Pesant S, Not F, Picheral M, Kandels-Lewis S, Le Bescot N, Gorsky G, Iudicone
637 D, Karsenti E, Speich S, Troublé R, Dimier C, Searson S, Acinas SG, Bork P,
638 Boss E, Bowler C, De Vargas C, Follows M, Gorsky G, Grimsley N, Hingamp P,
639 Iudicone D, Jaillon O, Kandels-Lewis S, Karp-Boss L, Karsenti E, Krzic U, Not F,
640 Ogata H, Pesant S, Raes J, Reynaud EG, Sardet C, Sieracki M, Speich S,
641 Stemmann L, Sullivan MB, Sunagawa S, Velayoudon D, Weissenbach J, Wincker
642 P. 2015. Open science resources for the discovery and analysis of Tara Oceans
643 data. *Open Sci Resour Discov Anal Tara Ocean data Sci Data* 2:150023.
- 644 19. Nurk S, Bankevich A, Antipov D. 2013. Assembling genomes and mini-
645 metagenomes from highly chimeric reads. *Res Comput Mol Biol* 158–170.
- 646 20. Edgar RC. 2007. PILER-CR: fast and accurate identification of CRISPR repeats.
647 *BMC Bioinformatics* 8:18.
- 648 21. Grissa I, Vergnaud G, Pourcel C. 2007. CRISPRFinder: a web tool to identify
649 clustered regularly interspaced short palindromic repeats. *Nucleic Acids Res*
650 35:W52-7.
- 651 22. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G,
652 Djahanschiri B, Zeller G, Mende DR, Alberti A, Cornejo-Castillo FM, Costea PI,
653 Cruaud C, D'Ovidio F, Engelen S, Ferrera I, Gasol JM, Guidi L, Hildebrand F,
654 Kokoszka F, Lepoivre C, Lima-Mendez G, Poulain J, Poulos BT, Royo-Llonch M,
655 Sarmiento H, Vieira-Silva S, Dimier C, Picheral M, Searson S, Kandels-Lewis S,

- 656 Bowler C, de Vargas C, Gorsky G, Grimsley N, Hingamp P, Iudicone D, Jaillon O,
657 Not F, Ogata H, Pesant S, Speich S, Stemmann L, Sullivan MB, Weissenbach J,
658 Wincker P, Karsenti E, Raes J, Acinas SG, Bork P. 2015. Ocean plankton.
659 Structure and function of the global ocean microbiome. *Science* 348:1261359.
- 660 23. Overbeek R, Begley T, Butler RM, Choudhuri J V., Chuang HY, Cohoon M, de
661 Crécy-Lagard V, Diaz N, Disz T, Edwards R, Fonstein M, Frank ED, Gerdes S,
662 Glass EM, Goesmann A, Hanson A, Iwata-Reuyl D, Jensen R, Jamshidi N,
663 Krause L, Kubal M, Larsen N, Linke B, McHardy AC, Meyer F, Neuweger H,
664 Olsen G, Olson R, Osterman A, Portnoy V, Pusch GD, Rodionov DA, Rülckert C,
665 Steiner J, Stevens R, Thiele I, Vassieva O, Ye Y, Zagnitko O, Vonstein V. 2005.
666 The subsystems approach to genome annotation and its use in the project to
667 annotate 1000 genomes. *Nucleic Acids Res* 33:5691–5702.
- 668 24. Seed KD, Lazinski DW, Calderwood SB, Camilli A. 2013. A bacteriophage
669 encodes its own CRISPR/Cas adaptive response to evade host innate immunity.
670 *Nature* 494:489–491.
- 671 25. Chenard C, Wirth JF, Suttle CA. 2016. Viruses infecting a freshwater filamentous
672 cyanobacterium (*Nostoc* sp.) encode a functional CRISPR array and a
673 proteobacterial DNA polymerase B. *MBio* 7:e00667-16.
- 674 26. Marraffini LA, Sontheimer EJ. 2010. Self versus non-self discrimination during
675 CRISPR RNA-directed immunity. *Nature* 463:568–71.
- 676 27. Garrett RA, Vestergaard G, Shah SA. 2011. Archaeal CRISPR-based immune
677 systems: Exchangeable functional modules. *Trends Microbiol* 19:549–556.
- 678 28. Heler R, Samai P, Modell JW, Weiner C, Goldberg GW, Bikard D, Marraffini LA.
679 2015. Cas9 specifies functional viral targets during CRISPR-Cas adaptation.
680 *Nature* 519:199–202.
- 681 29. Iranzo J, Lobkovsky AE, Wolf YI, Koonin E V. 2013. Evolutionary dynamics of the
682 prokaryotic adaptive immunity system CRISPR-Cas in an explicit ecological
683 context. *J Bacteriol* 195:3834–3844.
- 684 30. Sun CL, Thomas BC, Barrangou R, Banfield JF. 2016. Metagenomic
685 reconstructions of bacterial CRISPR loci constrain population histories. *ISME J*
686 10:858–870.
- 687 31. Ummat A, Bashir A. 2014. Resolving complex tandem repeats with long reads.
688 *Bioinformatics* 30:3491–3498.
- 689 32. Hara S, Koike I, Terauchi K, Kamiya H, Tanoue E. 1996. Abundance of viruses in
690 deep oceanic waters. *Mar Ecol Prog Ser* 145:269–277.
- 691 33. Westra ER, Van houte S, Oyesiku-Blakemore S, Makin B, Broniewski JM, Best A,
692 Bondy-Denomy J, Davidson A, Boots M, Buckling A. 2015. Parasite exposure
693 drives selective evolution of constitutive versus inducible defense. *Curr Biol*
694 25:1043–1049.
- 695 34. Morris RM, Rappé MS, Connon S a, Vergin KL, Siebold W a, Carlson C a,
696 Giovannoni SJ. 2002. SAR11 clade dominates ocean surface bacterioplankton
697 communities. *Nature* 420:806–810.
- 698 35. Zhao Y, Temperton B, Thrash JC, Schwalbach MS, Vergin KL, Landry ZC,
699 Ellisman M, Deerinck T, Sullivan MB, Giovannoni SJ. 2013. Abundant SAR11
700 viruses in the ocean. *Nature* 494:357–360.
- 701 36. Shah SA, Erdmann S, Mojica FJM, Garrett RA. 2013. Protospacer recognition

- 702 motifs: mixed identities and functional diversity. *RNA Biol* 10:891–9.
- 703 37. Krüger DH, Bickle T a. 1983. Bacteriophage survival: multiple mechanisms for
704 avoiding the deoxyribonucleic acid restriction systems of their hosts. *Microbiol*
705 *Rev* 47:345–360.
- 706 38. Minot S, Grunberg S, Wu GD, Lewis JD, Bushman FD. 2012. Hypervariable loci in
707 the human gut virome. *Proc Natl Acad Sci* 109:3962–3966.
- 708 39. Doulatov S, Hodes A, Dai L, Mandhana N, Zimmerly S, Miller JF, Liu M, Deora R,
709 Simons RW, Zimmerly S, Miller JF. 2004. Tropism switching in *Bordetella*
710 bacteriophage defines a family of diversity-generating retroelements. *Nature*
711 431:476–481.
- 712 40. Roux S, Hallam SJ, Woyke T, Sullivan MB. 2015. Viral dark matter and virus –
713 host interactions resolved from publicly available microbial genomes. *Elife* 4:1–20.
- 714 41. Emerson JB, Andrade K, Thomas BC, Norman A, Allen EE, Heidelberg KB,
715 Banfield JF. 2013. Virus-host and CRISPR dynamics in Archaea-dominated
716 hypersaline Lake Tyrrell, Victoria, Australia. *Archaea* 2013:370871.
- 717 42. Tschitschko B, Williams TJ, Allen MA, Páez-Espino D, Kyrpides N, Zhong L,
718 Raftery MJ, Cavicchioli R. 2015. Antarctic archaea–virus interactions:
719 metaproteome-led analysis of invasion, evasion and adaptation. *ISME J* 9:2094–
720 2107.
- 721 43. Paez-Espino D, Eloë-Fadrosch EA, Pavlopoulos GA, Thomas AD, Huntemann M,
722 Mikhailova N, Rubin E, Ivanova NN, Kyrpides NC. 2016. Uncovering Earth’s
723 virome. *Nature* 536:425–430.
- 724 44. Suzek BE, Huang H, McGarvey P, Mazumder R, Wu CH. 2007. UniRef:
725 Comprehensive and non-redundant UniProt reference clusters. *Bioinformatics*
726 23:1282–1288.
- 727 45. Angly FE, Willner D, Rohwer F, Hugenholtz P, Tyson GW. 2012. Grinder: a
728 versatile amplicon and shotgun sequence simulator. *Nucleic Acids Res* 40:e94.
- 729 46. Rho M, Wu YW, Tang H, Doak TG, Ye Y. 2012. Diverse CRISPRs evolving in
730 human microbiomes. *PLoS Genet* 8:e1002441.
- 731 47. Li W, Godzik A. 2006. Cd-hit: A fast program for clustering and comparing large
732 sets of protein or nucleotide sequences. *Bioinformatics* 22:1658–1659.
- 733 48. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat*
734 *Methods* 9:357–359.
- 735 49. Li H. 2011. A statistical framework for SNP calling, mutation discovery,
736 association mapping and population genetical parameter estimation from
737 sequencing data. *Bioinformatics* 27:2987–2993.
- 738 50. Wommack KE, Bhavsar J, Ravel J. 2008. Metagenomics: read length matters.
739 *Appl Env Microbiol* 74:1453–1463.
- 740 51. Skennerton CT, Imelfort M, Tyson GW. 2013. Crass: identification and
741 reconstruction of CRISPR from unassembled metagenomic data. *Nucleic Acids*
742 *Res* 41:e105.
- 743 52. Sorokin VA, Gelfand MS, Artamonova II, Artamonova II. 2010. Evolutionary
744 dynamics of clustered irregularly interspaced short palindromic repeat systems in
745 the ocean metagenome. *Appl Env Microbiol* 76:2136–2144.
- 746 53. Karsenti E, Acinas SG, Bork P, Bowler C, De Vargas C, Raes J, Sullivan M,
747 Arendt D, Benzoni F, Claverie J-M, Follows M, Gorsky G, Hingamp P, Iudicone D,

- 748 Jaillon O, Kandels-Lewis S, Krzic U, Not F, Ogata H, Pesant S, Reynaud EG,
749 Sardet C, Sieracki ME, Speich S, Velayoudon D, Weissenbach J, Wincker P.
750 2011. A Holistic Approach to Marine Eco-Systems Biology. PLoS Biol
751 9:e1001177.
- 752 54. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK,
753 Fierer N, Peña AG, Goodrich JK, Gordon JI, Huttley G a, Kelley ST, Knights D,
754 Koenig JE, Ley RE, Lozupone C a, Mcdonald D, Muegge BD, Pirrung M, Reeder
755 J, Sevinsky JR, Turnbaugh PJ, Walters W a, Widmann J, Yatsunencko T,
756 Zaneveld J, Knight R. 2010. QIIME allows analysis of high- throughput community
757 sequencing data. Nat Publ Gr 7:335–336.
- 758 55. Wommack KE, Bhavsar J, Polson SW, Chen J, Dumas M, Srinivasiah S, Furman
759 M, Jamindar S, Nasko DJ. 2012. VIROME: a standard operating procedure for
760 analysis of viral metagenome sequences. Stand Genomic Sci 6:427–439.
- 761 56. Pride DT, Salzman J, Relman DA. 2012. Comparisons of clustered regularly
762 interspaced short palindromic repeats and viromes in human saliva reveal
763 bacterial adaptations to salivary viruses. Environ Microbiol 14:2564–2576.
- 764 57. Manica A, Zebec Z, Steinkellner J, Schleper C. 2013. Unexpectedly broad target
765 recognition of the CRISPR-mediated virus defence system in the archaeon
766 *sulfolobus solfataricus*. Nucleic Acids Res 41:10509–10517.
- 767 58. Horvath P, Romero DA, Coute-Monvoisin A-C, Richards M, Deveau H, Moineau
768 S, Boyaval P, Fremaux C, Barrangou R. 2008. Diversity, Activity, and Evolution of
769 CRISPR Loci in *Streptococcus thermophilus*. J Bacteriol 190:1401–1412.
- 770 59. Noguchi H, Park J, Takagi T. 2006. MetaGene: Prokaryotic gene finding from
771 environmental genome shotgun sequences. Nucleic Acids Res 34:5623–5630.
- 772 60. Team RDC, R Development Core Team R. 2005. R: A Language and
773 Environment for Statistical Computing. R Found Stat Comput 1:409.
- 774 61. Hijmans RJ. 2014. Introduction to the " geosphere " package (Version 1 . 3-8) 1–
775 19.

776

777

778

779

780 **Acknowledgements**

781 This work was supported through grants to KEW and SWP from the National Science
782 Foundation (OCE-1148118, OIA-1736030 and DBI-1356374), the National Institutes for
783 Health (5R21AI109555-02) and the Gordon and Betty Moore Foundation (grant number
784 2732). Support from the University of Delaware Center for Bioinformatics and
785

786 Computational Biology Core Facility and use of the BIOMIX compute cluster was made
787 possible through funding from Delaware INBRE (NIGMS GM103446) and the Delaware
788 Biotechnology Institute.

789

790

791 **Authors' Contributions**

792 D.J.N, S.W.P., and K.E.W designed research; D.J.N. and R.M.M performed the
793 research; D.J.N. and J.D.B. wrote and modified software and D.J.N. and K.E.W. wrote
794 the paper. D.J.N., B.D.F., S.W.P., and K.E.W. revised the paper.

795

796 **Competing Financial Interests**

797 The authors declare no conflicts of interest in publishing this work.

798

799 **Material and Correspondence Information**

800 Address: Delaware Biotechnology Inst., 15 Innovation Way, Newark, Delaware 19711

801 (Tel): (302) 831-4362

802 (Fax): (302) 831-3447

803 (E-mail): wommack@dbi.udel.edu

804 **TABLES**

805 **Table 1:** Fifteen most abundant viroplankton ORFs containing viral gene target

806 sequences

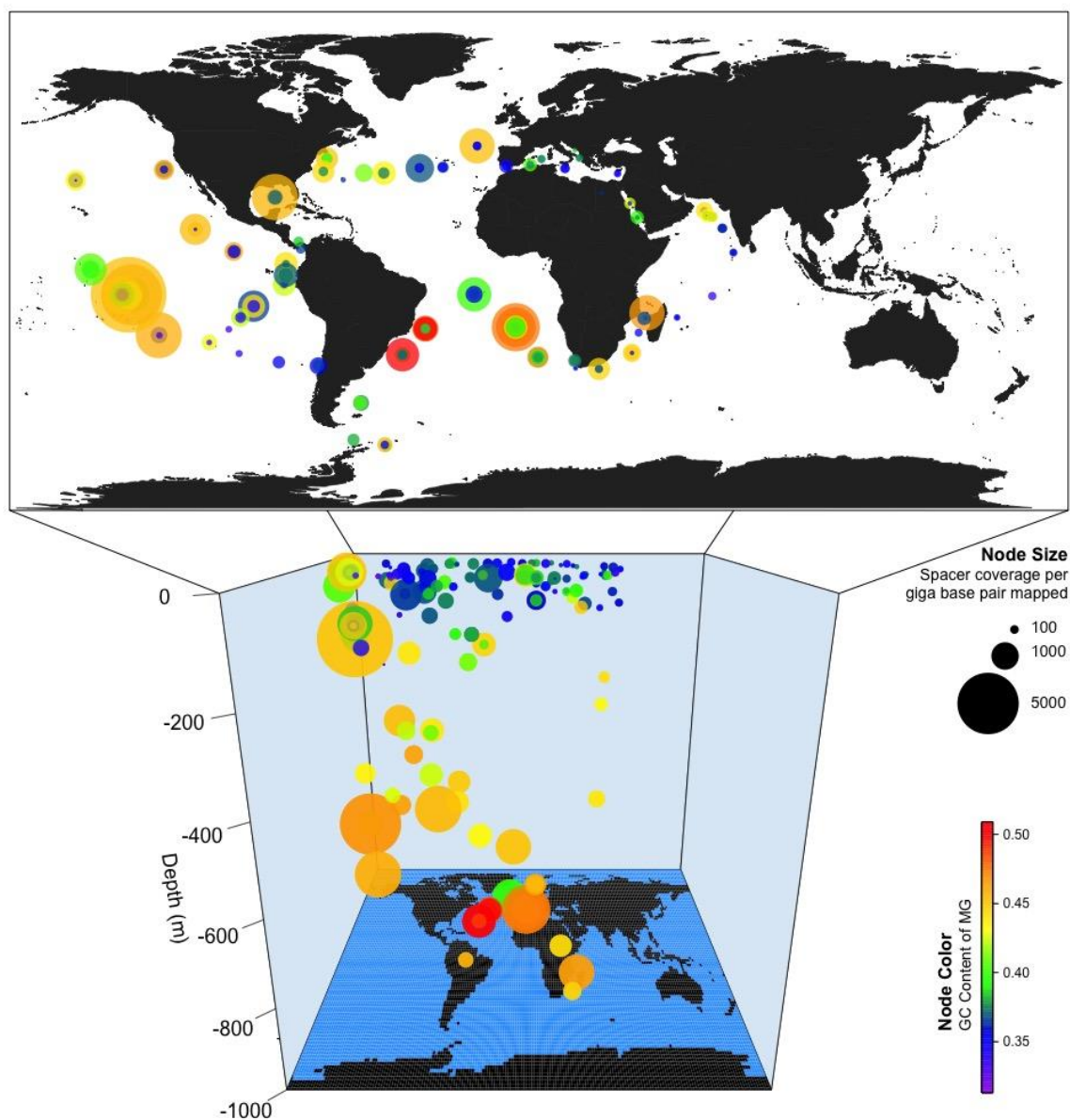
ORF Annotations	Actual Hits*	Expected Hits**	Fold (Act. / Exp.)
Phage protein	598	699	0.9
Phage terminase	102	48	2.1
Methyltransferase	67	14	4.7
Phage capsid protein	64	25	2.5
Phage tail protein	54	47	1.2
DNA polymerase	51	32	1.6
Phage-associated recombinase	37	12	3.0
Phage portal protein	35	24	1.5
ssDNA-binding protein	28	4	7.1
Phage DNA helicase	27	37	0.7
Reductase	27	23	1.2
Peptidase	23	18	1.3
Phage tail fiber	19	28	0.7
Phage tape measure protein	19	9	2.0
Glycotransferase	16	8	1.9
<i>[104 Other annotations]</i>	<i>272</i>	-	-
No hits	1920	2257	0.8

807 * Actual hits are the number of spacer hits

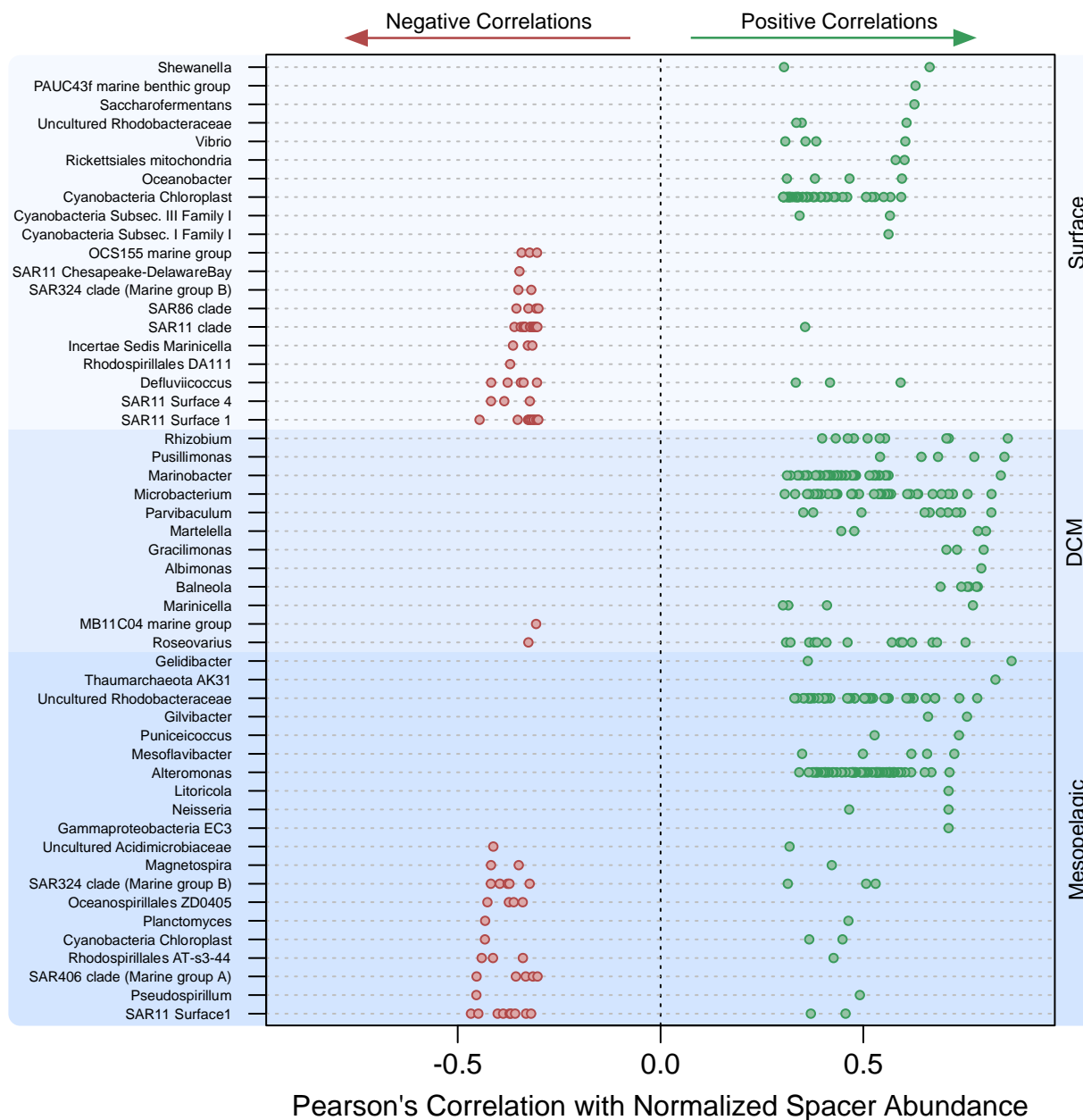
808 ** Expected hits were calculated based on the frequency of all aquatic viral ORFs being

809 assigned a given annotation by Phage SEED

810 **FIGURES**



811
812 **Fig. 1.** Normalized spacer abundance correlates with depth and GC content. Map of
813 spacers found by *Tara* Oceans sites. Node size represents the normalized abundance
814 of spacers at that *Tara* site (cumulative spacer coverage divided by mapped read
815 gigabases for that sample), node color represents the mean GC content of contigs at
816 that site.



817

818 **Fig. 2.** CRISPR abundance correlates with several taxonomic OTU's, with stronger
 819 positive correlations in deeper ocean zones. The top 10 positively and negatively
 820 correlating OTUs with respect to normalized spacer abundance, broken down by
 821 oceanic depth zone. Some taxa have several significant OTUs.

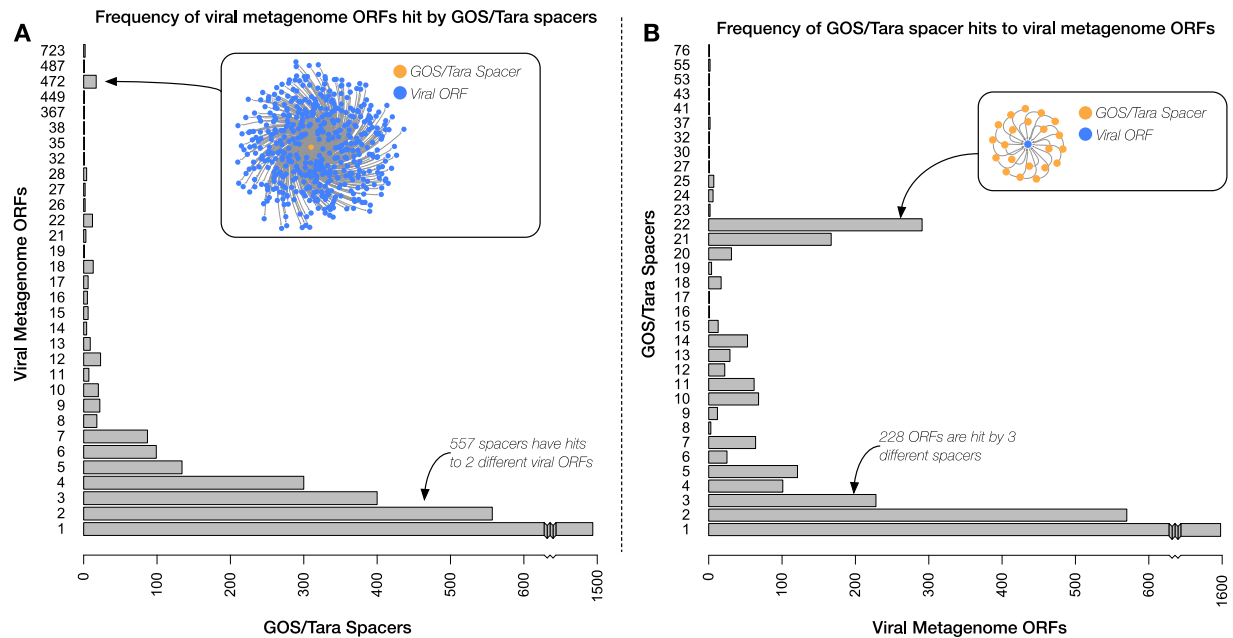
822

823

824

825

826



827

828 **Fig. 3.** CRISPR spacers aligned with viral gene targets in a many-to-many relationship.

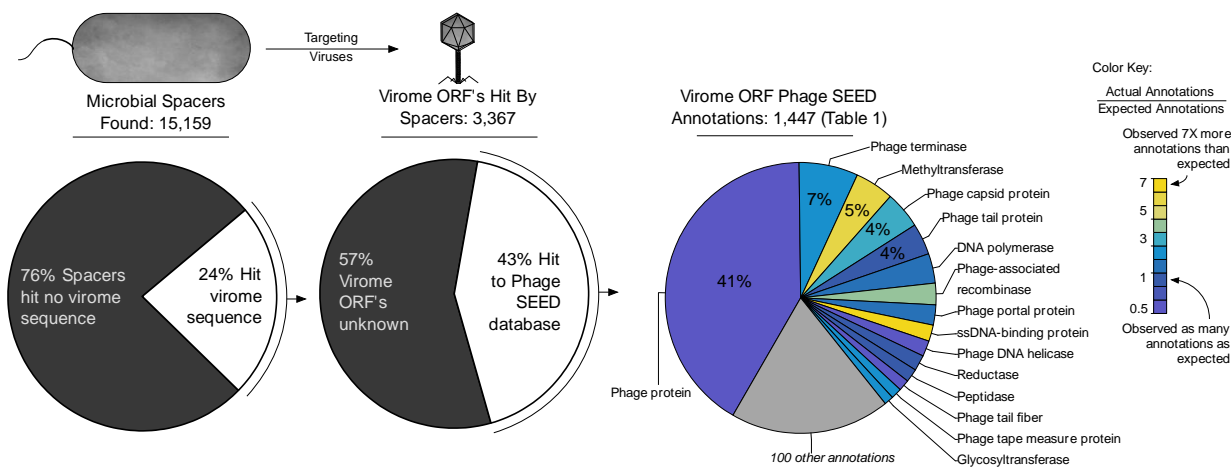
829 (A) Frequency of viral metagenome ORFs hit by GOS and *Tara* Oceans spacers with

830 inset network graph representing the 1 to 472 relationship. (B) Frequency of GOS and

831 *Tara* Oceans spacer hits to viral metagenome ORFs with inset network graph

832 representing the 22 to 1 relationship.

833
834
835
836
837

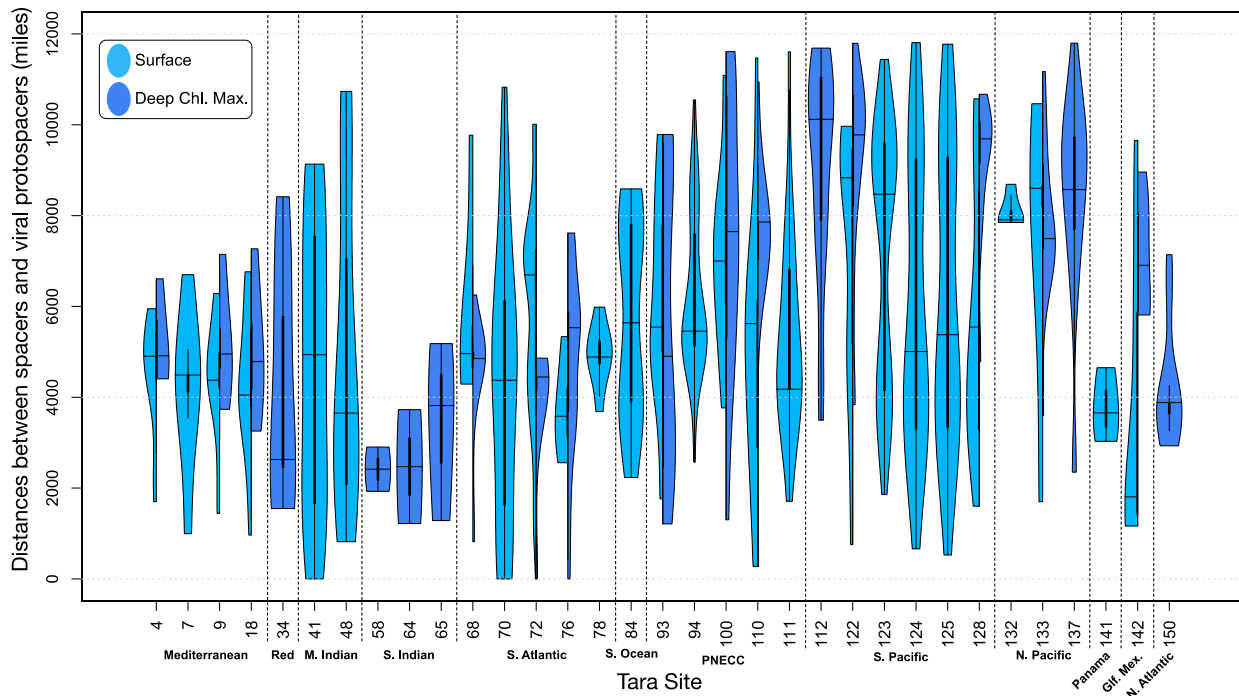


838

839 **Fig. 4.** Microbial spacers preferentially target specific viral genes. Nearly one quarter of
840 aquatic microbial spacers had a putative match to aquatic virome genes. The majority
841 of these genes obtained informative annotations (i.e. not “Phage protein”). Most genes
842 targeted by CRISPR spacers were annotated two-fold as often as expected, based on
843 the expected frequencies of aquatic virome gene annotations. Two gene annotations
844 that were seen less frequently than expected were DNA helicase and phage tail fiber.

845

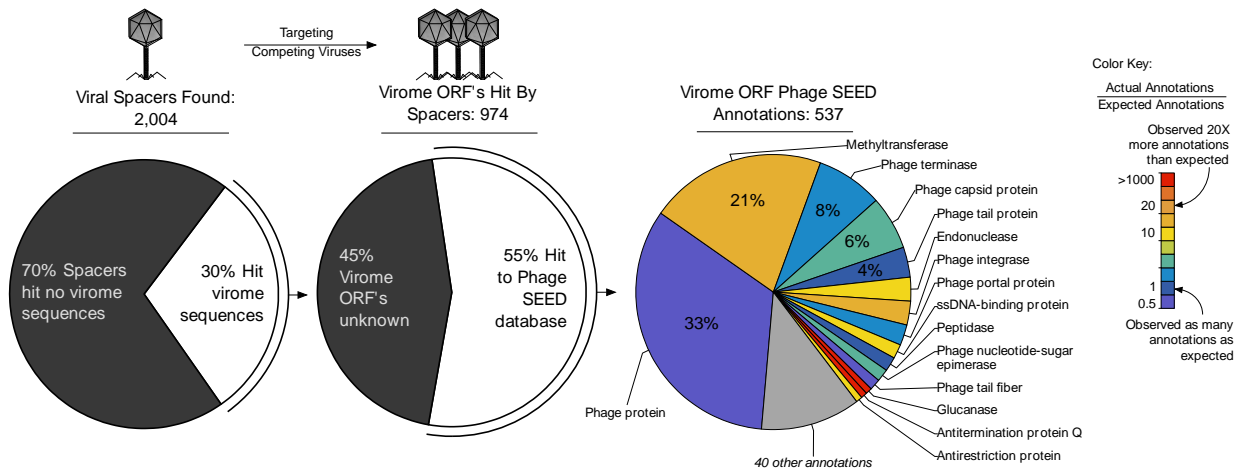
846



847

848 **Fig. 5.** CRISPR spacers are more likely to match viral gene targets from distant
849 viromes than paired viromes. Violin plots of the distances between *Tara* Oceans
850 spacers and the viromes they aligned with (light blue, surface sample; dark blue, deep
851 chlorophyll maximum sample). Line connections demonstrate sites with paired surface
852 and DCM samples. Sites are broadly split by geographic location (Mediterranean,
853 Mediterranean Sea; Red, Red Sea; M. Indian, Indian Monsoon Gyres; S. Indian, Indian
854 S. Subtropical Gyre; S. Atlantic, S. Atlantic Gyre; S. Ocean, Southern Ocean; PNECC,
855 Pacific North Equatorial Countercurrent; S. Pacific, South Pacific Ocean Gyre; N.
856 Pacific, North Pacific Ocean Gyre; Panama, near Panama; Gf. Mex., Gulf of Mexico; N.
857 Atlantic, North Atlantic Ocean Gyre).

858
859
860
861
862



863
864
865
866
867
868
869
870
871
872
873
874

Fig. 6. Over 2,000 CRISPR spacers were identified in the aquatic viral metagenomes and target methyltransferase more frequently than microbial spacers. Viral spacers are believed to assist the host in defending itself against competing viruses. Genes associated with temperate viruses (e.g. integrase, methyltransferase) are targeted more frequently by viral spacers than microbial spacers. Additionally, viral spacers targeted viral genes that were exceedingly rare in these aquatic dsDNA viromes, such as glucanase and antitermination protein Q, with many other genes being targeted >2X more often than expected. Again, phage tail fiber was targeted less frequently than expected.

875 **SUPPLEMENT**

876 **Figure S1:** The CASC Workflow. A) Preliminary search for CRISPR arrays and
877 identification of putative spacer arrays. B) Validation of putative spacers.

878 **Figure S2:** Nucleotide position histogram of CRISPR repeats from (A) CRISPR repeats
879 deemed “bona fide” by CASC, (B) all CRISPR repeats from CRISPR DB, and (C)
880 CRISPR repeats deemed Non-“bona fide” by CASC.

881 **Figure S3:** Alignments between spacers and viral ORFs were typically strong. (A)
882 Nearly 95% of HSPs had 3 or fewer mismatches in alignments of spacers to viral ORFs.
883 (B) Nearly 98% of HSPs had 1 or no gaps open in alignments between spaces and viral
884 ORFs.

885

886 **Additional File 1:** CRISPR spacers found in GOS datasets

887 **Additional File 2:** CRISPR spacers found in the Tara Oceans microbial metagenomes

888 **Additional File 3:** Summary of Tara Oceans viromes

889 **Additional File 4:** Summary of aquatic viromes collected from VIROME

890 (virome.dbi.udel.edu)

891 **Additional File 5:** Actual vs. expected number of annotations for candidate microbial-
892 viral protospacers

893 **Additional File 6:** Actual vs. expected number of annotations for candidate viral-viral
894 protospacers

895 **Table S1.** Bacterial genome sequences used in the construction of the mock metagenomes.

896

Organism	% GC	Genome Size (Mbp)	Spacers	Arrays
<i>Escherichia coli</i> str. K-12 substr. MG1655	51	4.6	18	2
<i>Streptococcus salivarius</i> JIM8777	40	2.2	32	1
<i>Neisseria meningitidis</i> 8013	51	2.3	25	1
<i>Yersinia pestis</i> A1122	48	4.5	16	3
<i>Chlorobium tepidum</i> TLS	57	2.1	62	2
<i>Chlamydia trachomatis</i> F/SW5	41	1.0	-	-
<i>Ruegeria pomeroyi</i> DSS-3	64	4.1	-	-
<i>Bacillus thuringiensis</i> HD-789	35	5.5	-	-
<i>Bordetella pertussis</i> CS	68	4.1	-	-
<i>Acetobacter pasteurianus</i> IFO 3283-01	53	2.9	-	-

897

898

899

900

901

902

903

904

905

906 **Table S2.** CRISPR finding tool performance Spacers found in the artificial 454 pyrosequencing metagenome using
 907 available CRISPR discovery tools.

908

Program	Spacers in Dataset	Spacers Detected	True Positives	False Positives	False Negatives	Sensitivity	Precision
CASC - Conservative	1930	981	802	179	1128	0.42	0.82
CASC - Liberal	1930	1623	1108	515	822	0.57	0.68
CRISPR Finder	1930	-	-	-	-	-	-
metaCRT	1930	2631	1225	1406	705	0.63	0.47
PILER-CR	1930	1483	1070	413	860	0.55	0.72

909
 910 **Table S3.** Spacers found in the artificial Illumina metagenome using available CRISPR discovery tools.

Program	Spacers in Dataset	Spacers Detected	Spacer Coverage in Dataset	Spacer Coverage Detected	True Positives	False Positives	False Negatives	Sensitivity	Precision
CASC - Conservative	153	153	42,349	41,095	153	0	0	1.00	1.00
CASC - Liberal	153	153	42,349	41,095	153	0	0	1.00	1.00
CRISPR Finder	153	216	42,349	62,418	153	63	0	1.00	0.71
metaCRT	153	365	42,349	96,503	153	212	0	1.00	0.42
PILER-CR	153	146	42,349	39,222	138	8	15	0.90	0.95

911
 912
 913