

1 Analysis of 19 Highly Conserved *Vibrio cholerae* 2 Bacteriophages Isolated from Environmental and 3 Patient Sources Over a Twelve-Year Period

4 Angus Angermeyer¹, Moon Moon Das², Durg Vijai Singh² and Kimberley D. Seed^{1,3,*}

5 ¹ Department of Plant and Microbial Biology, University of California Berkeley, Berkeley, CA, USA;
6 angermeyer@berkeley.edu (A.A.); kseed@berkeley.edu (K.D.S.)

7 ² Department of Infectious Disease Biology, Institute of Life Sciences, Nalco Square, Bhubaneswar 751023,
8 India; moon.scial@gmail.com (M.M.D.); durg.singh@gmail.com (D.V.S.)

9 ³ Chan Zuckerberg Biohub, San Francisco, CA, USA; kseed@berkeley.edu

10 * Correspondence: kseed@berkeley.edu; Tel.: +1-510-664-7711

11

12 **Abstract:** The *Vibrio cholerae* biotype 'El Tor' is responsible for all current epidemic and endemic
13 cholera outbreaks worldwide. These outbreaks are clonal and are hypothesized to originate from
14 the coastal areas near the Bay of Bengal where the lytic bacteriophage ICP1 specifically preys upon
15 these pathogenic outbreak strains. ICP1 has also been the dominant bacteriophage found in
16 cholera patient stool since 2001. However, little is known about its genomic differences between
17 ICP1 strains collected over time. Here we elucidate the pan-genome and phylogeny of ICP1 strains
18 by aligning, annotating and analyzing the genomes of 19 distinct isolates collected between 2001
19 and 2012. Our results reveal that ICP1 isolates are highly conserved and possess a large
20 core-genome as well as a smaller, somewhat flexible accessory-genome. Despite its overall
21 conservation, ICP1 strains have managed to acquire a number of unknown genes as well as a
22 CRISPR-Cas system, which is known to be critical for its ongoing struggle for co-evolutionary
23 dominance over its host. This study describes a foundation on which to construct future molecular
24 and bioinformatic studies of this *V. cholerae*-associated bacteriophages.

25 **Keywords:** Cholera; Bacteriophage; *Vibrio cholerae*; Pan-genome; Myovirus; Phylogenetics; *Vibrio*
26 phage

27

28 1. Introduction

29 *Vibrio cholerae* is a globally-distributed bacterium and the causal agent of the disease cholera, a
30 potentially severe intestinal illness that affects ~1-5 million people resulting in up to ~140,000 deaths
31 annually [1]. The current (seventh) pandemic is comprised of *V. cholerae* biotype 'El Tor'
32 (predominantly serotype O1) and is responsible for current endemic and epidemic disease [2].
33 Epidemic disease outbreaks sweep the globe periodically and have been traced back to a single
34 lineage that has emerged from the Bay of Bengal region in multiple waves over the last half-century
35 [3]. Despite the overall genetic heterogeneity of this lineage, in which individual outbreaks are
36 nearly always clonal [4], there is an abundance of subtle variation and horizontal transfer that has
37 been observed between outbreaks over time [5,6]. It is hypothesized that the Bay of Bengal serves as
38 a reservoir where El Tor strains circulate throughout the year exchanging genetic material and
39 undergoing ecological selection before infiltrating coastal communities. They are subsequently
40 transported by infected individuals to larger cities where they can be transmitted globally [5,7]. This
41 mechanism is thought to create a bottleneck for strains and result in the clonality of outbreaks.

42 Bacteriophage, viruses that uniquely infect bacteria, are extremely abundant in the
43 environment where they can outnumber their prokaryotic hosts by several orders of magnitude [8].
44 As such, bacteriophage play a key role in the evolution of their hosts through both selection and
45 phage-mediated lateral gene transfer [9]. These processes are likely to be very important to *V.*
46 *cholerae* strain evolution in the Bay of Bengal as well. Previous work has identified a *V. cholerae*

47 O1-specific [10] lytic *myoviridae* bacteriophage (ICP1) to be of particular interest in this system [11].
48 In Bangladesh, ICP1 has been found in water samples [12,13] and it has been identified as the
49 dominant phage in cholera patient stool samples since 2001 [11]. The persistence of this phage over
50 time indicates that *V. cholerae* has strategies to limit ICP1 predation, and that ICP1 can evolve to
51 overcome such defenses. Indeed, from this natural genetic laboratory, several complex and
52 surprising adaptations/acquisitions have occurred in the race for survival between *V. cholerae* and
53 ICP1. These include self-mobilizing chromosomal islands that can provide a rapid and efficient
54 response to ICP1 infection [14] and the first known example of a bacteriophage-encoded
55 CRISPR-Cas system [15]. Initial characterization of eight ICP1 isolates collected between 2001-2011
56 noted the relative low level of diversity and lack of major genomic rearrangements, deletions or
57 insertions [11] (with the exception of its remarkable CRISPR-Cas acquisition [15]). Here we build
58 upon the initial characterization of ICP1 to perform a comparative genomic analysis on 19
59 individual ICP1 isolates to reconstruct their phylogenetic relationships over time, identify the core
60 and accessory genomes, and infer possible gene function where possible.

61 *V. cholerae* is an organism that affects millions and appropriately, it is well-studied with
62 modern bioinformatic and sequencing tools. It is important that their concomitant bacteriophage
63 are similarly studied to help us better elucidate the important role that bacteriophage likely play in
64 the evolution of *V. cholerae* and the epidemiology of the ongoing cholera pandemic.
65

66 2. Materials and Methods

67 Nineteen ICP1 bacteriophage genomes were acquired from various sources (Table 1). Those
68 genomes not specifically described below were downloaded from the NCBI GenBank database.
69 Their metadata (if available) were used to inform our reporting of isolation source and isolation
70 year. In cases where this information differed from what was reported in a genome's original
71 publication we deferred to the GenBank database. Isolation year was used to standardize the ICP1
72 bacteriophage naming convention: ICP1_YEAR_X where "X" is sequentially assigned letter (A-Z)
73 based on the order in which genomes were named. The only exception is the original, 'ancestral'
74 ICP1 isolate, which is simply referred to as 'ICP1' [11].

75 Genomes not already available on GenBank include: ICP1_2006_E and ICP1_2011_A, which
76 were isolated and sequenced as described previously [15]. ICP1_2011_A was assembled using CLC
77 Genomics Workbench v10 (Qiagen, Redwood City, CA) and ICP1_2006_E was assembled with
78 IDBA-UD v1.1.3 [16] (default settings) after fastq read filtering with USEARCH v10.0.240 [17]
79 fastq_filter (-fastq_maxee_rate 0.001 -fastq_maxns 1 -fastq_truncqual 15 -fastq_maxee 0.25).
80 ICP1_2011_B was assembled from an existing diarrheal stool metagenomic sample [18] (SRA:
81 PRJEB9150; Run: ERR866578) using USEARCH filtering and IDBA-UD *de novo* assembly as above
82 (same parameters) to generate an incomplete ICP1 contig. That contig was then used as a reference
83 sequence to reassemble the same filtered reads using IDBA-Hybrid [16] with default settings.
84 Finally, ICP1_2012_A was isolated from a cholera patient stool sample collected in Silvassa, India
85 [19]. Genomic DNA libraries were sequenced using an Illumina Mi-Seq (Genotypic Technology,
86 Bangalore, India). Genome assembly was performed with CLC Genomics Workbench v10.

87 Whole-genome alignment was performed on all 19 genomes using progressiveMauve [20]
88 (build: Feb 25 2015) with default settings. The Mauve xmfa alignment output file was converted to
89 Phylip format using BioPython v1.71 [21] and a maximum likelihood, unrooted, phylogenetic tree
90 was constructed using PhyML v20120412 [22] while calculating bootstrap support (-s BEST
91 --rand_start --n_rand_starts 10 -b 100). A companion phylogeny based on the concatenated
92 core-genome (described below) was constructed using the same methods. Dendroscope3 v3.5.9 [23]
93 was used to generate a tanglegram joining both trees with ICP1 ancestral set as the outgroup. The
94 whole-genome alignment was also used to determine a consensus ICP1 sequence using CLC
95 Genomics Workbench v10, which all ICP1 genomes were visually mapped back onto using BRIG
96 v0.95 [24].

97 All extant CDS annotations, hereafter referred to as ‘Open Reading Frames’ (ORFs), in the ICP1
 98 ancestral GenBank file were blasted (BLASTn v2.6.0 [25]) against every other ICP1 genomic
 99 sequence to find homologous ORFs. Putative hits were considered homologs, and annotated as
 100 such, if the following conditions were met: the subject hit was a complete ORF (start and stop
 101 codons; codon table=11), $e\text{-value} \leq 1e-10$, $\%identity \geq 85$, and the two matched ORFs were within 10%
 102 sequence length of each other. After identifying existing homologs, *de novo* ORF prediction was
 103 performed on all genomes to find additional possible coding sequences using Prodigal v2.6.3 [26]
 104 with default settings and a confidence cutoff $\geq 95\%$. All newly identified putative-ORFs from each
 105 genome were then blasted against each other with the same conditions as above, grouped by
 106 homology and given an iterative numerical identifier based on their locations in the genome
 107 relative to the extant ORFs (i.e. ORF1, ORF2, ORF2.1, ORF2.2, ORF3, etc.). ORFs were then
 108 categorized into groups based on how many of the 19 genomes they were found in. If an ORF was
 109 present in all 19 genomes, it was considered part of the core-genome, otherwise, it was considered
 110 part of the accessory-genome. Core and accessory ORFs were also mapped to the BRIG alignment.
 111 The core-genome ORF protein sequences were concatenated by strain to create a core-genomic,
 112 pseudo-genome for each ICP1 strain and subjected to the phylogenetic analysis as above. The core
 113 and accessory ORFs were also queried against the NCBI’s Conserved Domain Database
 114 (<https://www.ncbi.nlm.nih.gov/cdd>) to determine if any possessed interesting domain homology
 115 that wasn’t detected by BLAST alone.

116
 117

Table 1. ICP1 bacteriophage strains.

Standardized Name	Previous Name	Isolation Year	Isolation Source	Genome Size (bp)	GenBank Accession	Genome Citation
ICP1	-	2001	Stool	125,956	HQ641347	Seed <i>et al.</i> 2011
ICP1_2001_A	-	2001	Stool	124,826	HQ641353	Seed <i>et al.</i> 2011
ICP1_2001_B	JSF1	2001	Water	126,082	KY883636	Naser <i>et al.</i> 2017
ICP1_2001_C	JSF2	2001	Water	126,082	KY883637	Naser <i>et al.</i> 2017
ICP1_2001_D	JSF4	2001*	Water*	124,261	KY065147	Naser <i>et al.</i> 2017
ICP1_2001_E	JSF5	2001*	Water	132,142	KY883634	Naser <i>et al.</i> 2017
ICP1_2001_F	JSF6	2001*	Water	133,685	KY883635	Naser <i>et al.</i> 2017
ICP1_2004_A	-	2004	Stool	128,083	HQ641354	Seed <i>et al.</i> 2011
ICP1_2005_A	-	2005	Stool	129,373	HQ641352	Seed <i>et al.</i> 2011
ICP1_2006_A	-	2006	Stool	123,104	HQ641351	Seed <i>et al.</i> 2011
ICP1_2006_B	-	2006	Stool	123,097	HQ641350	Seed <i>et al.</i> 2011
ICP1_2006_C	-	2006	Stool	124,497	HQ641349	Seed <i>et al.</i> 2011
ICP1_2006_D	-	2006	Stool	124,497	HQ641348	Seed <i>et al.</i> 2011
ICP1_2006_E	-	2006	Stool	128,298	TBD	This study
ICP1_2009_A	JSF13	2009	Water	128,814	KY883638	Naser <i>et al.</i> 2017
ICP1_2011_A	-	2011	Stool	126,861	TBD	This study
ICP1_2011_B	-	2011	Stool	125,128	TBD	This study
ICP1_2011_C	JSF14	2011	Water	125,096	KY883639	Naser <i>et al.</i> 2017
ICP1_2012_A	-	2012	Stool	121,418	TBD	This study

118

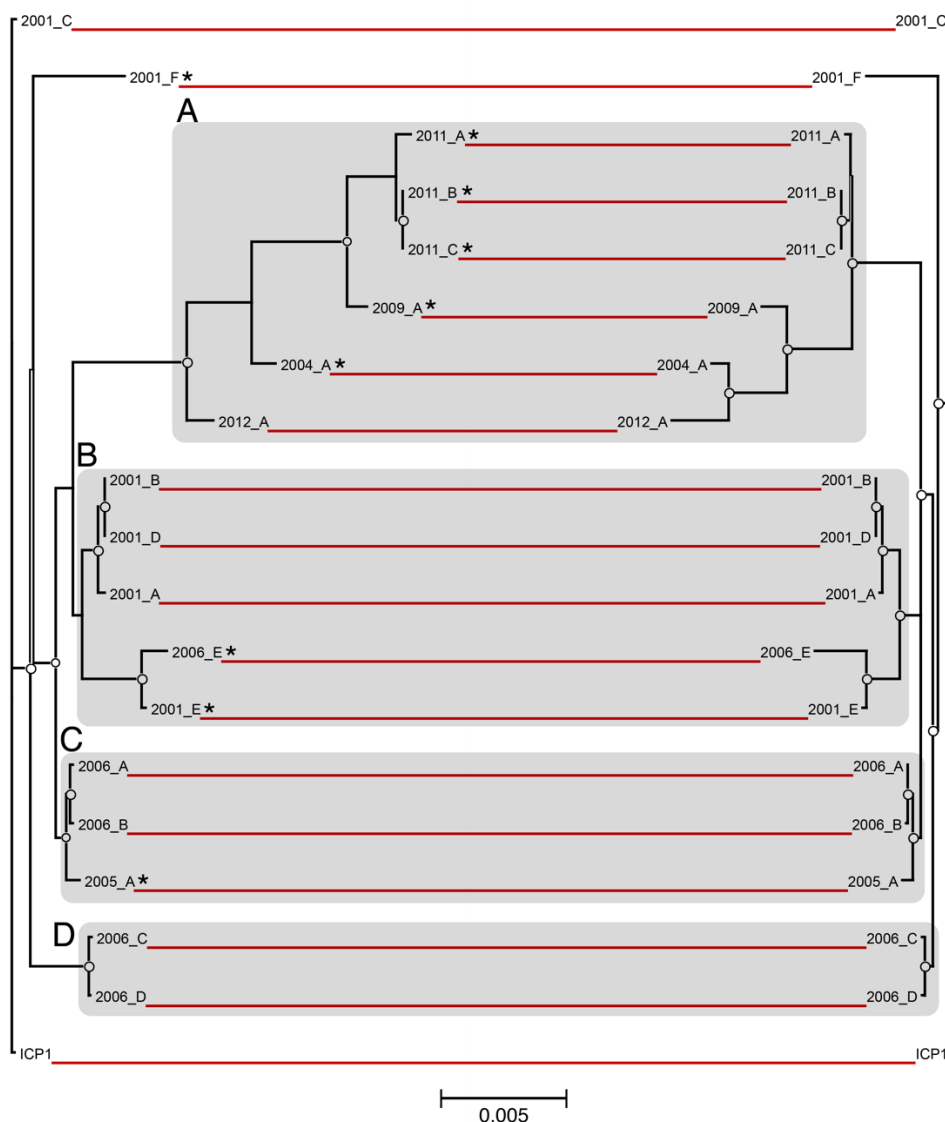
* GenBank metadata was reported in cases where it differed from the original publication.

119

120 3. Results

121 3.1. Genome Characteristics and Phylogeny

122 ICP1 genomes were analyzed from isolates collected over a 12-year period, from 2001 to 2012
123 (Table 1). The isolates were derived from both environmental water samples and patient stool
124 samples collected in Bangladesh (n=18) and India (ICP1_2012_A). Genome length was slightly
125 variable with an average of 126,384bp (stdev: 3008bp). The maximum likelihood phylogenetic
126 analysis grouped both the whole-genome and core-genome alignments into several general clusters
127 (Figure 1). Cluster A contains the five most recently isolated strains as well as ICP1_2004_A. This
128 cluster also contains five of the seven total CRISPR-positive strains in the dataset. Cluster B contains
129 four of the six 2001 isolates, ICP1_2006_E and two of the CRISPR-positive strains. Clusters C and D
130 contain isolates from intermediate years 2005 and 2006, with one CRISPR-positive represented in
131 cluster C. The CRISPR-positive strain ICP1_2001_F did not cluster with other isolates. ICP1
132 ancestral and ICP1_2001C also cluster closely together but are not specifically highlighted.



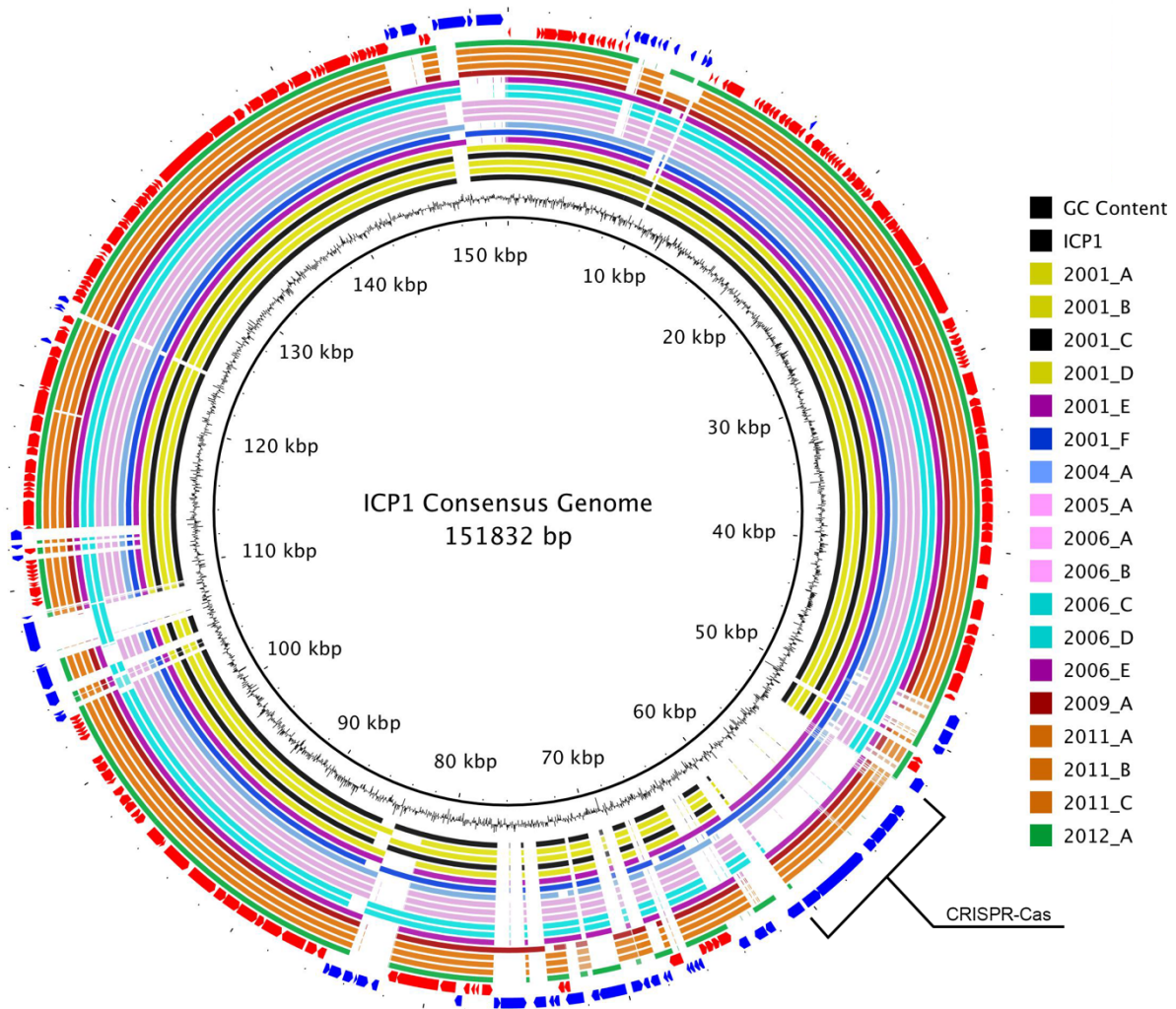
133
134
135
136
137
138
139
140

Figure 1. Phylogenetic comparison of ICP1 whole and core-genomes. Maximum likelihood phylogenetic trees (unrooted) were constructed based on a multiple whole-genome alignment of all 19 ICP1 phage sequences (left) as well as an alignment of the concatenated ORFs for each phage that comprise the core genome (right). The red lines connect identical leaves between trees to indicate relative phylogeny. The circles represent nodes with $\geq 90\%$ bootstrap support (n=100). The scale bar measures nucleotide substitutions per base pair. Distinct phylogenetic clusters are shaded grey and CRISPR strains are marked with '*'.

141 Overall, the topologies between whole-genome and core-genome trees were highly similar
142 with a few minor exceptions. Both share identical clustering and have almost no leaf-level
143 differences within clusters B, C and D. Cluster A showed an inversion of the phylogenetic
144 differences of the leaves between the two alignments. For instance, ICP1_2012_A has the most
145 divergent core-genome (from all other strains) but is the least divergent within its cluster when
146 the whole-genome was considered.

147 3.2. Genome Alignment Visualization

148 The consensus genomic sequence constructed from the whole-genome alignment was
149 151,832bp long and contained all of the coding and non-coding regions from each genome. It
150 was used as a reference to map the genomes and visualize the overall multiple-genome
151 alignment (Figure 2). Variable regions of insertions and deletions were visible as gaps in the
152 circular alignment. Similar to the whole-genome phylogeny, there was not a clear progression
153 of sequence divergence based on isolation chronology. No large regions of GC content
154 difference were observed.
155



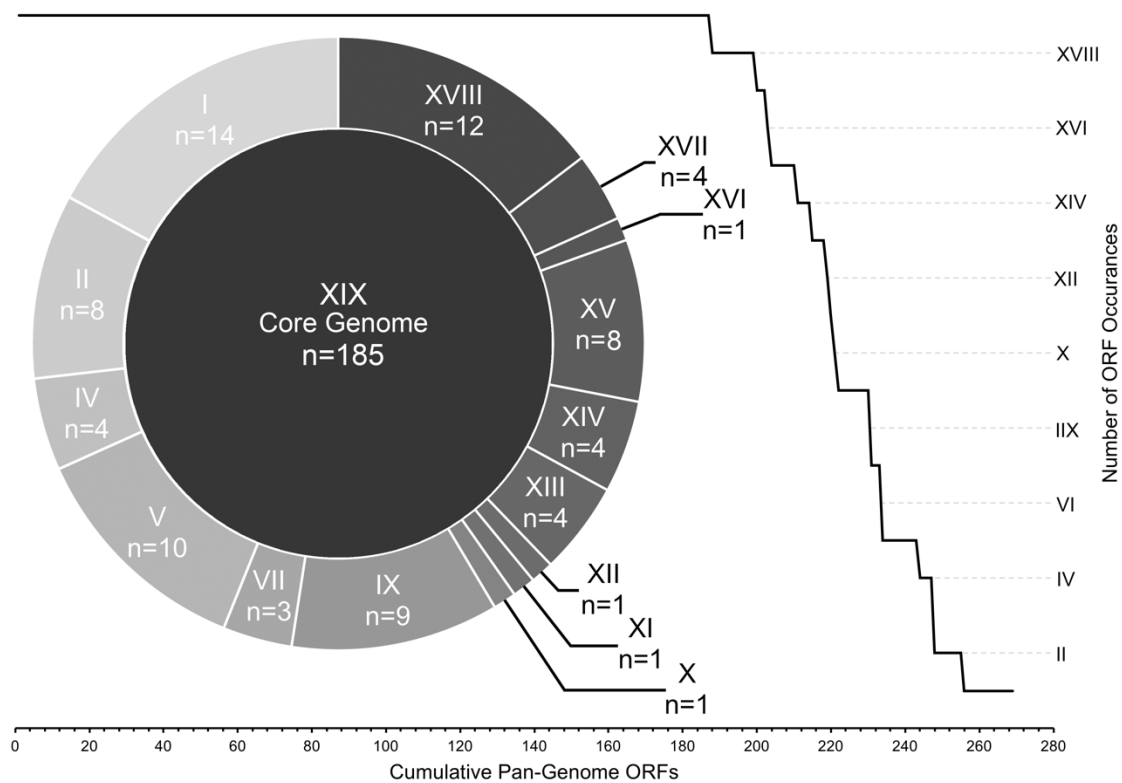
156
157 **Figure 2.** ICP1 pan-genome consensus alignment. A BLASTn-based whole genome alignment of all 19
158 ICP1 phage genomes using the MAUVE alignment consensus sequence as reference. The innermost
159 ring is the consensus sequence. The next ring represents the GC content for that region. The following
160 19 rings display the alignment for each genome and are colored by phylogenetic similarity as
161 determined by analysis of whole genomes (Figure 1 left). The second to last ring (red) represents the
162 core-genomic ORFs while the outermost ring (blue) is the accessory-genome ORFs. The CRISPR-Cas
163 insertion region is labeled.

164 3.3. ORF Annotation

165 Among all 19 genomes, a total of 269 distinct ORFs (based on homology cutoffs) were
166 identified. Of these, 230 were originally annotated in the ICP1 ancestral genome and 39 were
167 called by Prodigal. The number of ORFs per genome varied from 215 to 232 and demonstrated
168 a slight but significant inverse relationship with year of isolation, i.e. more recent strains had
169 fewer ORFs (Figure S1A). A weaker significant positive trend was observed between number
170 of ORFs and genome length (Figure S1B) which is likely due to the CRISPR-Cas insertions,
171 however there was not a statistically significant correlation between genome length and
172 isolation year.

173 3.3. Core and Accessory Genome Analysis

174 To estimate gene diversity among the genomes we divided the total pan-genomic ORF
175 complement into core and accessory groups. The core ORFs, those that were found in all 19
176 genomes, comprised ~70% of the total number of ORFs (185 out of 269) while the other 84 ORFs
177 were considered to be accessory (Figure 3).

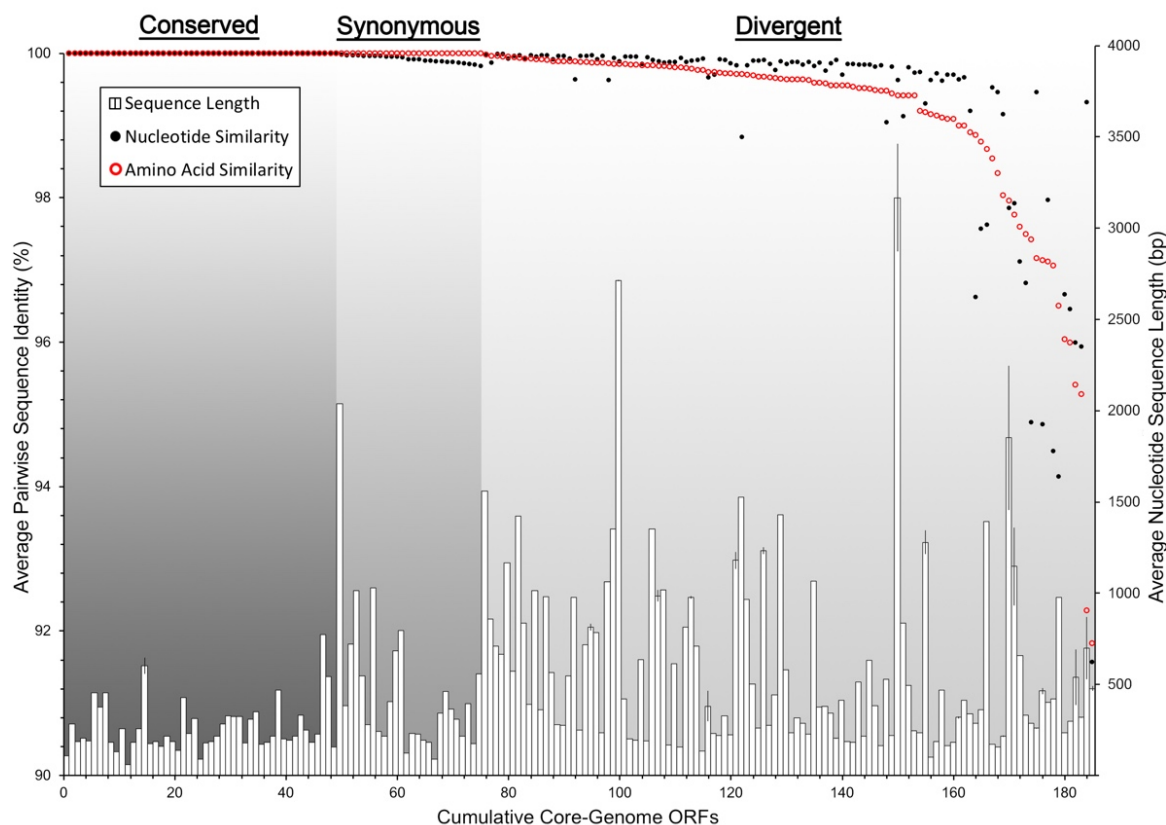


178

179 **Figure 3.** ICP1 pan-genome ORF allocation. All ICP1 ORFs arranged by number of genomes in
180 which they were detected. ORF bins are represented by roman numerals, i.e. 'XIX' ORFs were found
181 in all 19 genomes (core), 'X' ORFs were found in exactly 10 genomes, etc. The line graph
182 demonstrates the overall cumulative curve of core and accessory ORF prevalence. The doughnut
183 chart provides exact number of ORFs within each accessory bin.

184 Despite their occurrence across all genomes in this study, the core-genome ORFs were not all
185 equally conserved at the sequence level. By comparing average pairwise nucleotide and amino acid
186 sequence identity, we were able to resolve the core ORFs into three distinct groups: conserved-core,
187 synonymous-core and divergent-core (Figure 4). The conserved-core was comprised of 49 ORFs
188 that shared perfect sequence identity among all the genomic sequences. Correspondingly, they
189 shared identical amino acid sequence identity. The synonymous-core included 26 ORFs that

190 possessed a small amount of nucleotide diversity between genomes, but all mutations were silent,
191 and the amino acid primary structure was therefore identical. Finally, the divergent-core contained
192 the remaining 110 ORFs which were diverse in both nucleotide and amino acid pairwise identity.
193 These divergent ORFs encompassed a range pairwise identity similarities from 99.97% to 91.57% at
194 the nucleotide level and 99.98% to 91.83% for amino acid sequences (Table S1). The degree of
195 pairwise identity difference was not correlated with an ORF's sequence length (Figure 4).



196

197 **Figure 4.** ICP1 core-genome ORF divergence. All ICP1 core-genome ORFs arranged by average
198 pairwise similarity (171 pairwise comparisons) of both DNA nucleotide sequence alignments (black
199 dots) and amino acid residue alignments (red circles). The histogram shows average nucleotide
200 sequence length and standard deviation among the 19 sequences per ORF. The graph is divided
201 into three sections by types of ORF similarity: Conserved (identical nucleotide and amino acid
202 sequences), Synonymous (identical amino acid sequences, but silent nucleotide mutations) and
203 Divergent (dissimilarity in both nucleotide and amino acid sequences).

204 The accessory ORFs were distributed among all 19 genomes in a complex manner (Table S2)
205 and grouped into 15 levels of occurrence. These ranged from occurrences in 18 genomes (n=12) to
206 singletons (n=14) (Figure 3) with several patterns of accessory ORF co-occurrence (Table S2). The
207 most obvious example was the CRISPR-Cas locus which was previously known to reside in 8 of
208 these genomes [13,15] and confirmed through our annotation methods. Other examples of
209 sequential ORF co-occurrence included a locus containing ORFs 115, 116, 117, 117.1 and 118 which
210 were found in the same five genomes, ORFs 160, 162, 163 in a different set of five overlapping
211 genomes and ORFs 222, 223, 225 in 15 genomes.

212 3.3. Conserved Functional Domains

213 The majority of ORFs in both the core and accessory genomes were classified as
214 hypothetical due to a lack of an informative BLAST identification. Only 18 of the core ORFs
215 (9.7%) currently had a predicted function, two in the conserved-core, three in the

216 synonymous-core and 13 in the divergent-core (Table S1). The accessory core contains 11 ORFs
217 with a putative or known function (Table S2). The NCBI conserved domain search identified 18
218 additional core-genome ORFs and 11 accessory-genome ORFs that shared at least partial amino
219 acid sequence homology to known functional domains.

220 4. Discussion

221 ICP1 is a *V. cholerae*-infecting bacteriophage that appears to be prevalent throughout the Bay of
222 Bengal's coastal areas and is transferred readily alongside *V. cholerae* into humans during cholera
223 outbreaks in the region [11]. It is becoming increasingly theorized that this region is the source of
224 most if not all global cholera outbreaks [3,4,27], and therefore a better understanding of this
225 co-evolving predator is an important area of ongoing cholera research. In this study we have
226 performed a comprehensive phylogenetic analysis on all available, well-sequenced ICP1 isolates to
227 elucidate their genetic divergence over time and provide a platform on which to develop future
228 ICP1-related bioinformatic analyses.

229 We have found that the genomes of ICP1 are surprisingly well-conserved between all isolates
230 over the twelve-year period in which they were isolated. This is demonstrated in the whole-genome
231 phylogeny which, while resolvable into distinct phylogenetic clusters, still only represents a
232 maximum variation of approximately 1 nucleotide substitution per 100 base pairs between the most
233 divergent isolates (Figure 1), many of which are likely silent or non-coding. A high degree of
234 genomic conservation is also indicated by the relatively large core-genome shared between all
235 isolates (Figure 3). This conservation is not only surprising due to the amount of time the
236 core-genome has remained stable, but also due to the complex conditions that likely exist in the
237 coastal and ocean environments where ICP1 is competing against a host population that is almost
238 certainly more diverse than clonal outbreak strains. In at least one other study that examined
239 multiple strains of a single marine bacteriophage, Far-T4, it was shown that sequence variability
240 among strains was at least an order of magnitude greater than for ICP1 [28]. In contrast, a study
241 from a less variable environment found that strains of a *Y. enterocolitica*-infecting podovirus were
242 more highly conserved than ICP1 [29]. However, it must be considered that these environmental
243 pressures may actually be what drives this bacteriophage to maintain a large core-genome, of
244 which the vast majority of ORFs are hypothetical. This large gene complement may be providing
245 the flexibility needed to compete in the Bay of Bengal. Currently we can only speculate on the exact
246 reason for the conservation of ICP1's genome, but expanding the repertoire of well-sequenced
247 genomes will help to put a finer constraint on the core-genome and perhaps identify genetic targets
248 for future study.

249 Despite the stable conservation of the core-genome, ICP1 also possesses a diverse collection of
250 accessory genes that have been integrated into several locations around the pan-genome (Figure 2).
251 These accessory ORFs appear to be both single acquisitions as well as integrations of larger loci
252 (Table S2). The most notable of the latter is the previously described acquisition of a CRISPR-Cas
253 system which is used as a weapon in the arms race between ICP1 and *V. cholerae* [14,15]. This may
254 suggest that other mechanisms of molecular warfare would be likely targets for acquisition, but if
255 so, the conserved domain analyses failed to reveal mechanisms already known. Interestingly,
256 though, the trend has been for genomic ORF counts to diminish over time culminating in the latest
257 isolate from India in 2012 which possessed the fewest overall number of ORFs, i.e. the smallest
258 accessory-genome (Figure S1). This could be due to shedding of outdated or detrimental genes,
259 possibly because the acquisition of a CRISPR-Cas system provided enough flexibility to obviate the
260 need for other mechanisms. And although ICP1_2012_A is CRISPR-Cas negative, the other five
261 most recent isolates are positive. It is also possible that this is a regional difference though, and
262 more sequenced genomes will help to determine if this trend is chronological or geographical or
263 spurious.

264 Querying the ORFs against the NCBI's conserved domain database returned several hits having
265 to do with replication, nucleotide metabolism, recombination, endonuclease activity and a few
266 other basic functions (Table S1, S2). However, what may be most telling is that 80% of all ORFs do

267 not possess homology to any conserved domain indicating that there is a great deal left to learn
268 about the interaction between ICP1 and *V. cholerae*. It should also be noted that the annotation of
269 ORFs necessarily requires certain assumptions to be made about similarity cutoffs and thresholds
270 and as such these ORF calls are best viewed as estimates. However, we were conservative in our
271 methods and are confident that the a very large proportion of the calls are accurate.

272 As we advance our understanding of how cholera spreads globally, it will be important to also
273 continue tracking ICP1's phylogeny and genetic composition so that we may develop a better
274 understanding of its co-evolution with *V. cholerae* and attempt to disentangle the complex molecular
275 and ecological interactions that may play an important role in defining cholera outbreaks.

276 **Supplementary Materials:** The following are available online at www.mdpi.com/xxx/s1, Figure S1: ORF
277 occurrence trends by year and genome length, Table S1: Core-genome ORFs pairwise similarity and CDD hits,
278 Table S2: Accessory-genome ORFs ICP1 strain matrix.

279 **Author Contributions:** Angus Angermeyer and Kimberley Seed conceived of the study; Angus Angermeyer
280 analyzed the data; Moon Moon Das and Durg Vijai Singh contributed sequencing data and reagents. Angus
281 Angermeyer and Kimberley Seed wrote the paper.

282 **Funding:** This research was funded by the National Institute of Allergy and Infectious Diseases grant number
283 R01AI127652 and the Chan Zuckerberg Biohub.

284 **Acknowledgments:** We wish to thank Zach Barth for his helpful assistance and **Genotypic Technology,**
285 **Bangalore, India** for the sequencing of ICP1_2012_A.

286 **Conflicts of Interest:** The authors declare no conflict of interest. The funding sponsors had no role in the
287 design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, and
288 in the decision to publish the results.

289 References

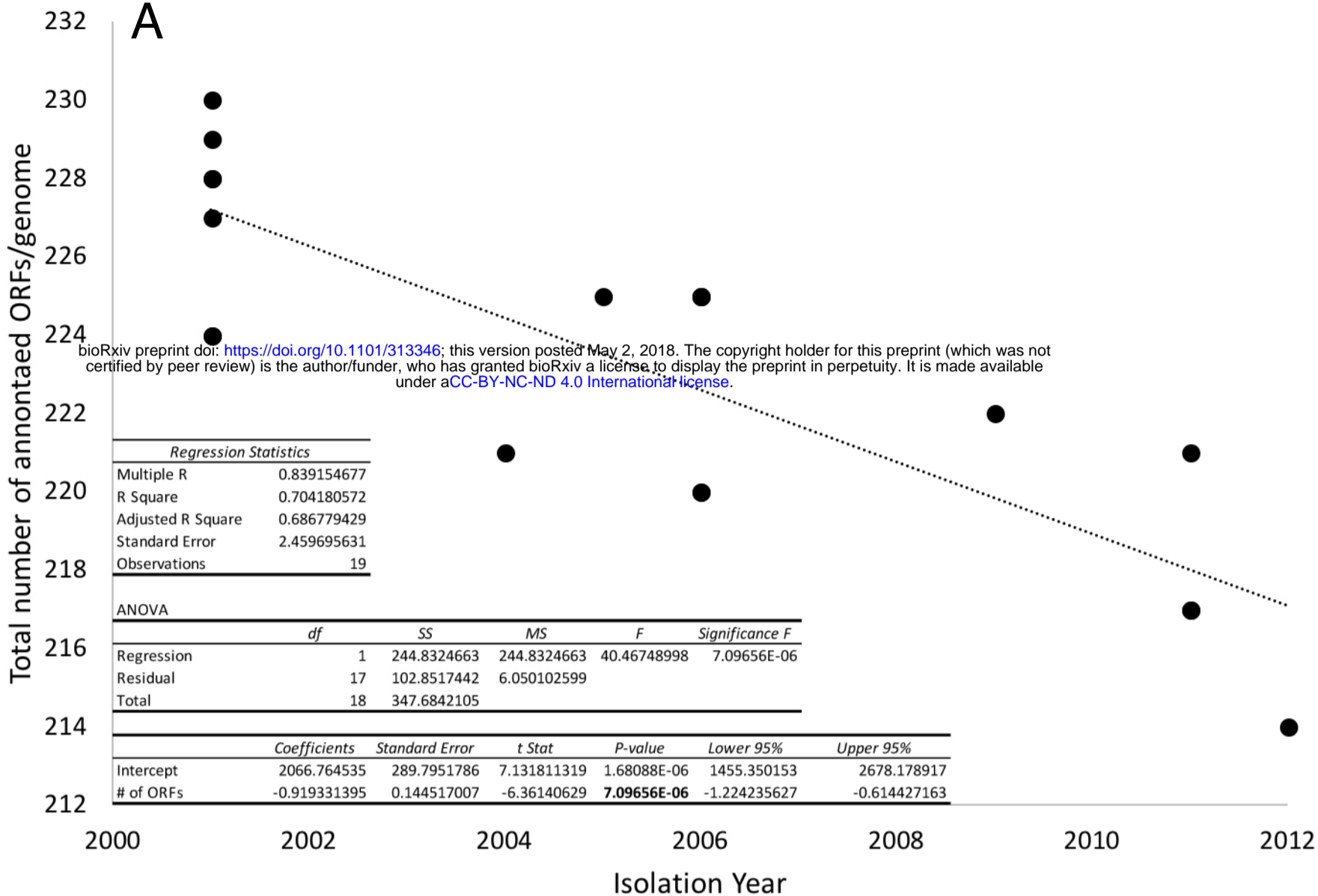
- 290 1. Ali, M.; Nelson, A. R.; Lopez, A. L.; Sack, D. A. Updated global burden of cholera in endemic
291 countries. *PLoS Negl Trop Dis* **2015**, *9*, e0003832.
- 292 2. Safa, A.; Nair, G. B.; Kong, R. Y. C. Evolution of new variants of *Vibrio cholerae* O1. *Trends*
293 *Microbiol.* **2010**, *18*, 46–54.
- 294 3. Mutreja, A.; Kim, D. W.; Thomson, N. R.; Connor, T. R.; Lee, J. H.; Kariuki, S.; Croucher, N. J.;
295 Choi, S. Y.; Harris, S. R.; Lebens, M.; Niyogi, S. K.; Kim, E. J.; Ramamurthy, T.; Chun, J.; Wood, J.
296 L. N.; Clemens, J. D.; Czerkinsky, C.; Nair, G. B.; Holmgren, J.; Parkhill, J.; Dougan, G. Evidence
297 for several waves of global transmission in the seventh cholera pandemic. *Nature* **2011**, *477*,
298 462–465.
- 299 4. Moore, S.; Thomson, N.; Mutreja, A.; Piarroux, R. Widespread epidemic cholera caused by a
300 restricted subset of *Vibrio cholerae* clones. *Clin. Microbiol. Infect.* **2014**, *20*, 373–379.
- 301 5. Weill, F.-X.; Domman, D.; Njamkepo, E.; Tarr, C.; Rauzier, J.; Fawal, N.; Keddy, K. H.; Salje, H.;
302 Moore, S.; Mukhopadhyay, A. K.; Bercion, R.; Luquero, F. J.; Ngandjio, A.; Dosso, M.;
303 Monakhova, E.; Garin, B.; Bouchier, C.; Pazzani, C.; Mutreja, A.; Grunow, R.; Sidikou, F.; Bonte,
304 L.; Breurec, S.; Damian, M.; Njanpop-Lafourcade, B.-M.; Sapriel, G.; Page, A.-L.; Hamze, M.;
305 Henkens, M.; Chowdhury, G.; Mengel, M.; Koeck, J.-L.; Fournier, J.-M.; Dougan, G.; Grimont, P.
306 A. D.; Parkhill, J.; Holt, K. E.; Piarroux, R.; Ramamurthy, T.; Quilici, M.-L.; Thomson, N. R.
307 Genomic history of the seventh pandemic of cholera in Africa. *Science* **2017**, *358*, 785–789.
- 308 6. Cho, Y.-J.; Yi, H.; Lee, J. H.; Kim, D. W.; Chun, J. Genomic evolution of *Vibrio cholerae*. *Curr Opin*
309 *Microbiol* **2010**, *13*, 646–651.
- 310 7. Domman, D.; Quilici, M.-L.; Dorman, M. J.; Njamkepo, E.; Mutreja, A.; Mather, A. E.; Delgado, G.;
311 Morales-Espinosa, R.; Grimont, P. A. D.; Lizárraga-Partida, M. L.; Bouchier, C.; Aanensen, D. M.;

- 312 Kuri-Morales, P.; Tarr, C. L.; Dougan, G.; Parkhill, J.; Campos, J.; Cravioto, A.; Weill, F.-X.;
313 Thomson, N. R. Integrated view of *Vibrio cholerae* in the Americas. *Science* **2017**, *358*, 789–793.
- 314 8. Parikka, K. J.; Le Romancer, M.; Wauters, N.; Jacquet, S. Deciphering the virus-to-prokaryote ratio
315 (VPR): insights into virus-host relationships in a variety of ecosystems. *Biol Rev Camb Philos Soc*
316 **2017**, *92*, 1081–1100.
- 317 9. Ofir, G.; Sorek, R. Contemporary Phage Biology: From Classic Models to New Insights. *Cell* **2018**,
318 *172*, 1260–1270.
- 319 10. Seed, K. D.; Faruque, S. M.; Mekalanos, J. J.; Calderwood, S. B.; Qadri, F.; Camilli, A. Phase
320 Variable O Antigen Biosynthetic Genes Control Expression of the Major Protective Antigen and
321 Bacteriophage Receptor in *Vibrio cholerae* O1. *PLoS Pathog* **2012**, *8*, e1002917–13.
- 322 11. Seed, K. D.; Bodi, K. L.; Kropinski, A. M.; Ackermann, H.-W.; Calderwood, S. B.; Qadri, F.;
323 Camilli, A. Evidence of a dominant lineage of *Vibrio cholerae*-specific lytic bacteriophages shed
324 by cholera patients over a 10-year period in Dhaka, Bangladesh. *MBio* **2011**, *2*, e00334–10.
- 325 12. Faruque, S. M.; Naser, I. B.; Islam, M. J.; Faruque, A. S. G.; Ghosh, A. N.; Nair, G. B.; Sack, D. A.;
326 Mekalanos, J. J. Seasonal epidemics of cholera inversely correlate with the prevalence of
327 environmental cholera phages. *Proc Natl Acad Sci USA* **2005**, *102*, 1702–1707.
- 328 13. Naser, I. B.; Hoque, M. M.; Nahid, M. A.; Tareq, T. M.; Rocky, M. K.; Faruque, S. M. Analysis of
329 the CRISPR-Cas system in bacteriophages active on epidemic strains of *Vibrio cholerae* in
330 Bangladesh. *Scientific Reports* **2017**, *7*, 14880.
- 331 14. O'Hara, B. J.; Barth, Z. K.; McKitterick, A. C.; Seed, K. D. A highly specific phage defense system
332 is a conserved feature of the *Vibrio cholerae* mobilome. *PLoS Genet.* **2017**, *13*, e1006838.
- 333 15. Seed, K. D.; Lazinski, D. W.; Calderwood, S. B.; Camilli, A. A bacteriophage encodes its own
334 CRISPR/Cas adaptive response to evade host innate immunity. *Nature* **2013**, *494*, 489–491.
- 335 16. Peng, Y.; Leung, H. C. M.; Yiu, S. M.; Chin, F. Y. L. IDBA-UD: a de novo assembler for single-cell
336 and metagenomic sequencing data with highly uneven depth. *Bioinformatics* **2012**, *28*, 1420–1428.
- 337 17. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **2010**,
338 *26*, 2460–2461.
- 339 18. David, L. A.; Weil, A.; Ryan, E. T.; Calderwood, S. B.; Harris, J. B.; Chowdhury, F.; Begum, Y.;
340 Qadri, F.; LaRocque, R. C.; Turnbaugh, P. J. Gut microbial succession follows acute secretory
341 diarrhea in humans. *MBio* **2015**, *6*, e00381–15.
- 342 19. Das, M. M.; Bhotra, T.; Zala, D.; Singh, D. V. Phenotypic and genetic characteristics of *Vibrio*
343 *cholerae* O1 carrying Haitian ctxB and attributes of classical and El Tor biotypes isolated from
344 Silvassa, India. *Journal of Medical Microbiology* **2016**, *65*, 720–728.
- 345 20. Darling, A. E.; Mau, B.; Perna, N. T. progressiveMauve: multiple genome alignment with gene
346 gain, loss and rearrangement. *PLoS ONE* **2010**, *5*, e11147.
- 347 21. Cock, P. J. A.; Antao, T.; Chang, J. T.; Chapman, B. A.; Cox, C. J.; Dalke, A.; Friedberg, I.;
348 Hamelryck, T.; Kauff, F.; Wilczynski, B.; de Hoon, M. J. L. Biopython: freely available Python
349 tools for computational molecular biology and bioinformatics. *Bioinformatics* **2009**, *25*, 1422–1423.
- 350 22. Guindon, S.; Dufayard, J.-F.; Lefort, V.; Anisimova, M.; Hordijk, W.; Gascuel, O. New algorithms
351 and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML
352 3.0. *Syst. Biol.* **2010**, *59*, 307–321.
- 353 23. Huson, D. H.; Scornavacca, C. Dendroscope 3: an interactive tool for rooted phylogenetic trees
354 and networks. *Syst. Biol.* **2012**, *61*, 1061–1067.

- 355 24. Alikhan, N.-F.; Petty, N. K.; Ben Zakour, N. L.; Beatson, S. A. BLAST Ring Image Generator
356 (BRIG): simple prokaryote genome comparisons. *BMC Genomics* **2011**, *12*, 402.
- 357 25. Camacho, C.; Coulouris, G.; Avagyan, V.; Ma, N.; Papadopoulos, J.; Bealer, K.; Madden, T. L.
358 BLAST+: architecture and applications. *BMC Bioinformatics* **2009**, *10*, 421.
- 359 26. Hyatt, D.; Chen, G.-L.; LoCascio, P. F.; Land, M. L.; Larimer, F. W.; Hauser, L. J. Prodigal:
360 prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* **2010**,
361 *11*, 119.
- 362 27. Hu, D.; Liu, B.; Feng, L.; Ding, P.; Guo, X.; Wang, M.; Cao, B.; Reeves, P. R.; Wang, L. Origins of
363 the current seventh cholera pandemic. *Proc Natl Acad Sci USA* **2016**, *113*, E7730–E7739.
- 364 28. Roux, S.; Enault, F.; Ravet, V.; Pereira, O.; Sullivan, M. B. Genomic characteristics and
365 environmental distributions of the uncultivated Far-T4 phages. *Front Microbiol* **2015**, *6*, 199.
- 366 29. Salem, M.; Skurnik, M. Genomic Characterization of Sixteen *Yersinia enterocolitica*-Infecting
367 Podoviruses of Pig Origin. *Viruses* **2018**, *10*.
- 368
- 369

Number of ORFs in Genome vs. Isolation Year

A



Number of ORFs in Genome vs. Genome Length

B

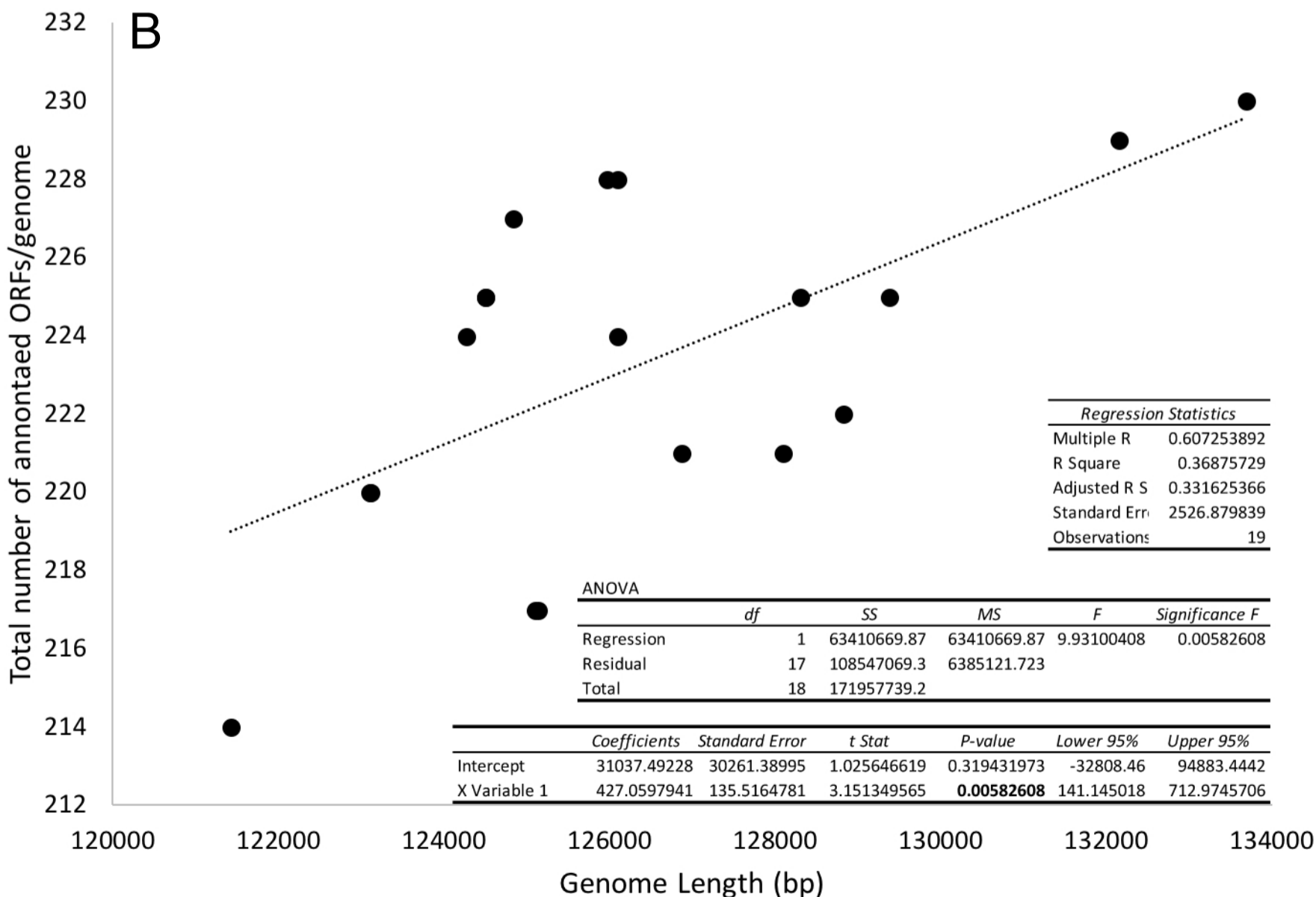


Figure S1: ORF linear regression statistics. ANOVA analyses of the linear regressions between (A) number of ORFs/genome vs genome isolation year and (B) ORFs/genome vs. genome length.

Table S1: Core-genome Information

	ORF Name	Nucleotide Similarity	Amino Acid Similarity	nucl_stddev	aa_stddev	Known Function	CDD Domain name	E-Value	PSSM-ID
Conserved Core	ORF1	100.00	100.00	0.00	0.00	-	-	-	-
	ORF10	100.00	100.00	0.00	0.00	-	-	-	-
	ORF109	100.00	100.00	0.00	0.00	-	-	-	-
	ORF12	100.00	100.00	0.00	0.00	-	-	-	-
	ORF13	100.00	100.00	0.00	0.00	-	-	-	-
	ORF134	100.00	100.00	0.00	0.00	-	-	-	-
	ORF135	100.00	100.00	0.00	0.00	-	-	-	-
	ORF138	100.00	100.00	0.00	0.00	-	-	-	-
	ORF142	100.00	100.00	0.00	0.00	-	-	-	-
	ORF154	100.00	100.00	0.00	0.00	-	-	-	-
	ORF158	100.00	100.00	0.00	0.00	-	-	-	-
	ORF164	100.00	100.00	0.00	0.00	-	-	-	-
	ORF167	100.00	100.00	0.00	0.00	-	-	-	-
	ORF168	100.00	100.00	0.00	0.00	-	-	-	-
	ORF171	100.00	100.00	0.00	0.00	-	-	-	-
	ORF175	100.00	100.00	0.00	0.00	-	HAD_like superfamily	1.87E-04	328728
	ORF180	100.00	100.00	0.00	0.00	-	-	-	-
	ORF195	100.00	100.00	0.00	0.00	-	-	-	-
	ORF206	100.00	100.00	0.00	0.00	-	-	-	-
	ORF214	100.00	100.00	0.00	0.00	-	-	-	-
	ORF216	100.00	100.00	0.00	0.00	-	-	-	-
	ORF218	100.00	100.00	0.00	0.00	-	-	-	-
	ORF226	100.00	100.00	0.00	0.00	-	-	-	-
	ORF227	100.00	100.00	0.00	0.00	-	-	-	-
	ORF26	100.00	100.00	0.00	0.00	-	-	-	-
	ORF29	100.00	100.00	0.00	0.00	-	-	-	-
	ORF3	100.00	100.00	0.00	0.00	-	-	-	-
	ORF32	100.00	100.00	0.00	0.00	-	-	-	-
	ORF33	100.00	100.00	0.00	0.00	-	-	-	-
	ORF34	100.00	100.00	0.00	0.00	-	-	-	-
	ORF39	100.00	100.00	0.00	0.00	-	-	-	-
	ORF41	100.00	100.00	0.00	0.00	-	DUF3696 superfamily	6.88E-03	331172
	ORF42	100.00	100.00	0.00	0.00	-	-	-	-
	ORF47	100.00	100.00	0.00	0.00	-	-	-	-
	ORF48	100.00	100.00	0.00	0.00	-	-	-	-
	ORF50	100.00	100.00	0.00	0.00	-	-	-	-
	ORF52	100.00	100.00	0.00	0.00	-	DUF1778 superfamily	5.16E-06	321696
	ORF54	100.00	100.00	0.00	0.00	-	-	-	-
	ORF55	100.00	100.00	0.00	0.00	ribonuclease H	RNase_Hi prokaryote_like	1.19E-56	260010
	ORF56	100.00	100.00	0.00	0.00	-	RNase_H_like superfamily	2.12E-14	326352
	ORF60	100.00	100.00	0.00	0.00	-	-	-	-
	ORF61	100.00	100.00	0.00	0.00	-	-	-	-
	ORF62	100.00	100.00	0.00	0.00	-	-	-	-
	ORF64	100.00	100.00	0.00	0.00	-	-	-	-
	ORF65	100.00	100.00	0.00	0.00	-	-	-	-
	ORF66	100.00	100.00	0.00	0.00	-	Peptidases_S8_S53 superfamily	5.78E-03	324584
	ORF75	100.00	100.00	0.00	0.00	putative baseplate assembly protein	Phage_base_V superfamily	9.89E-03	327437
	ORF82	100.00	100.00	0.00	0.00	-	-	-	-
ORF9	100.00	100.00	0.00	0.00	-	-	-	-	
Synonymous Core	ORF57	99.98	100.00	0.02	0.00	putative primase/helicase	RecA-like_NTPases superfamily	7.80E-28	333705
	ORF166	99.97	100.00	0.08	0.00	-	-	-	-
	ORF80	99.97	100.00	0.06	0.00	HNH homing endonuclease	HNHc superfamily	7.94E-05	320750
	ORF189	99.97	100.00	0.05	0.00	-	DnaQ_like_exo superfamily	1.39E-05	324557
	ORF137	99.96	100.00	0.07	0.00	-	-	-	-
	ORF213	99.96	100.00	0.11	0.00	-	-	-	-
	ORF81	99.96	100.00	0.05	0.00	-	Macolilin superfamily	3.48E-04	313022
	ORF43	99.96	100.00	0.13	0.00	-	-	-	-
	ORF45	99.95	100.00	0.14	0.00	-	-	-	-
	ORF74	99.95	100.00	0.10	0.00	-	-	-	-
	ORF4	99.95	100.00	0.07	0.00	-	-	-	-
	ORF112	99.94	100.00	0.06	0.00	putative DNA-binding protein Roi	Phage_pRha superfamily	2.26E-08	324635
	ORF191	99.91	100.00	0.25	0.00	-	-	-	-
	ORF46	99.91	100.00	0.17	0.00	-	-	-	-
	ORF30	99.91	100.00	0.18	0.00	-	-	-	-
	ORF44	99.90	100.00	0.20	0.00	-	-	-	-
	ORF210	99.89	100.00	0.22	0.00	-	-	-	-
	ORF153	99.88	100.00	0.34	0.00	-	-	-	-
	ORF139	99.88	100.00	0.14	0.00	-	-	-	-
	ORF83	99.88	100.00	0.12	0.00	-	-	-	-
	ORF7	99.88	100.00	0.14	0.00	-	-	-	-
	ORF194	99.87	100.00	0.16	0.00	-	-	-	-
	ORF2	99.86	100.00	0.22	0.00	-	-	-	-
	ORF181	99.85	100.00	0.22	0.00	-	-	-	-
	ORF15	99.84	100.00	0.25	0.00	-	-	-	-
	ORF126	99.82	100.00	0.36	0.00	-	-	-	-
	ORF176	99.97	99.98	0.03	0.06	putative exodeoxyribonuclease	-	-	-
	ORF173	99.87	99.96	0.09	0.11	recombination-associated protein RdgC	RdgC superfamily	5.57E-21	321354
	ORF170	99.99	99.96	0.04	0.13	-	-	-	-
	ORF179	99.98	99.95	0.05	0.14	-	-	-	-
	ORF68	99.92	99.95	0.11	0.11	-	-	-	-
	ORF49	99.96	99.94	0.07	0.16	-	-	-	-
	ORF73	99.96	99.94	0.04	0.10	putative baseplate component	Baseplate_J superfamily	1.29E-05	321435
	ORF196	99.91	99.92	0.12	0.22	putative adenine methyltransferase	Dam	2.70E-32	223415
	ORF131	99.97	99.92	0.08	0.24	-	-	-	-
	ORF76	99.95	99.92	0.07	0.13	-	-	-	-
	ORF71	99.97	99.91	0.08	0.26	-	-	-	-
	ORF211	99.97	99.91	0.05	0.15	ClpP ATP-dependent protease subunit	crotonase-like superfamily	6.26E-19	329030
	ORF120	99.91	99.89	0.10	0.22	-	-	-	-
	ORF63	99.96	99.89	0.11	0.33	-	-	-	-
	ORF132	99.96	99.88	0.11	0.34	-	-	-	-
	ORF177	99.92	99.88	0.13	0.23	-	-	-	-
	ORF79	99.63	99.88	0.55	0.16	-	-	-	-
	ORF140	99.96	99.87	0.12	0.37	-	DUF2130 superfamily	5.56E-03	331406
	ORF207	99.96	99.87	0.07	0.20	putative Gp5 baseplate hub subunit and tail lysozyme	NLPC_P60 superfamily	3.15E-35	328779
	ORF59	99.97	99.87	0.08	0.38	-	-	-	-
	ORF72	99.91	99.87	0.08	0.18	-	-	-	-

Divergent Core

ORF188	99.96	99.86	0.13	0.40	-	-	-	-
ORF124	99.63	99.86	0.46	0.16	-	-	-	-
ORF204	99.93	99.85	0.06	0.14	ribonucleoside diphosphate reductase, beta chain	Ferritin_like superfamily	1.00E-52	320867
ORF58	99.89	99.85	0.09	0.12	DNA polymerase	DNA_pol_A superfamily	5.25E-33	322025
ORF169	99.95	99.85	0.10	0.30	-	-	-	-
ORF202	99.95	99.84	0.15	0.46	-	-	-	-
ORF144	99.95	99.84	0.16	0.48	-	-	-	-
ORF36	99.85	99.83	0.14	0.23	-	-	-	-
ORF11	99.95	99.83	0.16	0.49	-	-	-	-
ORF215	99.90	99.83	0.08	0.14	DNA ligase	CDC9 superfamily	1.57E-99	330238
ORF28	99.88	99.82	0.23	0.37	-	HNHc	1.08E-05	238038
ORF122	99.86	99.82	0.12	0.22	putative major head protein	Phage_cap_E superfamily	6.03E-24	309113
ORF155	99.87	99.81	0.37	0.56	-	-	-	-
ORF78	99.88	99.81	0.13	0.26	-	-	-	-
ORF151	99.93	99.80	0.19	0.59	-	-	-	-
ORF208	99.88	99.80	0.12	0.22	PhoH family protein	P-loop_NTPase superfamily	4.93E-31	328724
ORF53	99.89	99.78	0.16	0.31	-	-	-	-
ORF172	99.91	99.77	0.14	0.33	-	-	-	-
ORF25	99.92	99.76	0.23	0.70	-	-	-	-
ORF192	99.66	99.74	0.45	0.49	-	-	-	-
ORF145	99.70	99.74	0.46	0.53	-	-	-	-
ORF27	99.91	99.73	0.18	0.55	-	-	-	-
ORF8	99.91	99.72	0.15	0.44	-	-	-	-
ORF6	99.86	99.72	0.31	0.83	-	-	-	-
ORF130	99.83	99.71	0.12	0.26	-	-	-	-
ORF127	98.84	99.71	2.24	0.47	-	-	-	-
ORF193	99.83	99.70	0.13	0.25	putative thymidylate synthase	Thy1 superfamily	1.56E-04	332234
ORF152	99.90	99.69	0.10	0.30	-	-	-	-
ORF220	99.89	99.67	0.17	0.52	-	-	-	-
ORF121	99.90	99.67	0.07	0.21	-	TSorf172	8.07E-05	313714
ORF143	99.85	99.67	0.26	0.53	-	-	-	-
ORF133	99.77	99.66	0.18	0.40	-	-	-	-
ORF128	99.88	99.65	0.11	0.34	terminase large subunit	Terminase_6 superfamily	7.55E-28	321850
ORF40	99.85	99.64	0.14	0.34	-	-	-	-
ORF217	99.88	99.64	0.19	0.58	-	-	-	-
ORF184	99.88	99.63	0.16	0.50	-	-	-	-
ORF187	99.84	99.63	0.23	0.50	-	-	-	-
ORF150	99.88	99.63	0.20	0.60	-	-	-	-
ORF70	99.82	99.59	0.09	0.30	-	-	-	-
ORF31	99.86	99.59	0.17	0.51	-	-	-	-
ORF185	99.75	99.58	0.21	0.40	-	-	-	-
ORF51	99.85	99.55	0.17	0.52	-	-	-	-
ORF212	99.90	99.55	0.19	0.72	-	-	-	-
ORF77	99.70	99.55	0.63	0.64	-	-	-	-
ORF156	99.85	99.55	0.24	0.72	-	-	-	-
ORF136	99.85	99.53	0.25	0.75	-	-	-	-
ORF200	99.84	99.52	0.19	0.58	-	-	-	-
ORF186	99.84	99.51	0.22	0.66	-	-	-	-
ORF209	99.84	99.50	0.14	0.43	-	-	-	-
ORF123	99.81	99.48	0.24	0.57	-	-	-	-
ORF108	99.83	99.48	0.27	0.83	-	-	-	-
ORF107	99.04	99.48	1.41	0.64	-	-	-	-
ORF199	99.82	99.44	0.22	0.67	-	-	-	-
ORF203	99.63	99.42	0.83	1.27	-	Ribonuc_red_1gC superfamily	3.16E-87	332162
ORF125	99.13	99.42	0.94	0.46	-	Peptidase_578_2 superfamily	4.65E-10	317012
ORF205	99.81	99.42	0.19	0.56	-	-	-	-
ORF219	99.72	99.41	0.26	0.70	-	-	-	-
ORF37	99.74	99.19	0.26	0.78	-	-	-	-
ORF165	99.30	99.18	1.32	1.43	-	-	-	-
ORF129	99.62	99.15	0.51	1.36	-	-	-	-
ORF141	99.72	99.14	0.27	0.82	-	-	-	-
ORF198	99.62	99.10	0.29	0.63	-	-	-	-
ORF201	99.70	99.09	0.30	0.93	-	-	-	-
ORF67	99.70	99.08	0.32	0.97	-	-	-	-
ORF110	99.63	99.00	0.63	1.60	-	-	-	-
ORF69	99.67	98.99	0.22	0.68	-	-	-	-
ORF94	99.20	98.90	1.98	2.19	-	-	-	-
ORF103	96.62	98.87	3.49	0.87	-	-	-	-
ORF93	97.56	98.77	3.12	1.64	-	-	-	-
ORF84	97.62	98.67	3.12	1.99	-	DUF3383 superfamily	4.95E-09	314693
ORF189.1	99.52	98.54	0.58	1.76	-	-	-	-
ORF197	99.46	98.34	0.42	1.30	-	-	-	-
ORF157	99.15	98.03	0.94	1.79	-	-	-	-
ORF113	97.85	97.96	2.06	2.60	-	NRDD superfamily	5.49E-149	330954
ORF174	97.92	97.76	8.27	8.40	-	TSorf172	1.47E-04	313714
ORF91	97.11	97.59	2.66	2.23	-	-	-	-
ORF92	96.81	97.49	3.76	2.96	-	-	-	-
ORF161	94.88	97.41	6.15	3.11	-	-	-	-
ORF35	99.45	97.16	1.10	5.70	-	-	-	-
ORF87	94.86	97.13	4.65	2.77	-	-	-	-
ORF119	97.97	97.12	2.17	2.95	-	-	-	-
ORF86	94.49	97.05	5.22	2.98	-	-	-	-
ORF85	94.13	96.50	7.18	4.47	-	-	-	-
ORF159	96.66	96.04	3.20	3.89	-	BF2867_like_C	2.18E-03	240526
ORF104	96.45	96.00	5.71	6.21	-	-	-	-
ORF221	95.99	95.40	5.24	5.88	-	-	-	-
ORF97	95.94	95.27	6.13	6.91	-	PP-binding superfamily	1.14E-07	324546
ORF5	99.32	92.28	0.85	10.28	-	GIY-YIG_SF superfamily	5.49E-03	326551
ORF114	91.57	91.83	6.46	6.32	-	Radical_SAM superfamily	1.93E-74	327492

Table S1: Core-genome information. The data for each ORF represented in Figure 4 is listed in columns 2 and 3. Columns 4 and 5 contain the standard deviations for those pairwise similarity values. The other columns contain information about putative gene function and conserved domain homology.

