

1 Nature Microbiology. Letter

2

3 **Origins of pandemic *Vibrio cholerae* from environmental gene pools**

4 B. Jesse Shapiro<sup>1</sup>, Inès Levade<sup>1#</sup>, Gabriela Kovacikova<sup>2#</sup>, Ronald K. Taylor<sup>2‡</sup>, Salvador

5 Almagro-Moreno<sup>2,3\*</sup>

6

7 <sup>1</sup>Department of Biological Sciences, University of Montreal, Montreal, Quebec, Canada.

8 <sup>2</sup>Department of Microbiology and Immunology, Geisel School of Medicine at Dartmouth,

9 Hanover, New Hampshire, USA. <sup>3</sup>Burnett School of Biomedical Sciences, College of Medicine,

10 University of Central Florida, Orlando, Florida, USA.

11

12 \*For correspondence: samoreno@ucf.edu

13 #These authors contributed equally

14 ‡Deceased

15 **Abstract**

16 Some microbes can transition from an environmental lifestyle to a pathogenic one<sup>1-3</sup>. This  
17 ecological switch typically occurs through the acquisition of horizontally acquired virulence  
18 genes<sup>4,5</sup>. However, the genomic features that must be present in a population prior to the  
19 acquisition of virulence genes and emergence of pathogenic clones remain unknown. We  
20 hypothesized that virulence adaptive polymorphisms (VAPs) circulate in environmental  
21 populations and are required for this transition. We developed a comparative genomic  
22 framework for identifying VAPs, using *Vibrio cholerae* as a model. We then characterized  
23 several environmental VAP alleles to show that, while some of them reduced the ability of  
24 clinical strains to colonize a mammalian host, other alleles conferred efficient host colonization.  
25 These results show that VAPs are present in environmental bacterial populations prior to the  
26 emergence of virulent clones. We propose a scenario in which VAPs circulate in the  
27 environment, they become selected and enriched under certain ecological conditions, and finally  
28 a genomic background containing several VAPs acquires virulence factors that allows for its  
29 emergence as a pathogenic clone.

## 30 **Main text**

31 Numerous bacterial pathogens have emerged from environmental populations<sup>1-3,6</sup>. These  
32 virulent clones evolve through the acquisition of toxins and host colonization factors<sup>4,5</sup>. Given  
33 that the genes encoding these factors can often spread widely by horizontal gene transfer, it is  
34 surprising that only a limited number of pathogenic clones have emerged from any particular  
35 bacterial species. As a model of how environmental gene pools give rise to pandemic clones,  
36 we used *Vibrio cholerae*, a genetically diverse group of aquatic bacteria that include a confined  
37 phylogenetic group, the “pandemic genome” group (PG), that can cause the severe diarrheal  
38 disease cholera in humans<sup>7,8</sup>. Seven pandemics of cholera have been recorded to date, all  
39 caused by the PG group. The current pandemic is caused by strains of the El Tor biotype, and  
40 has spread across the globe in several waves of transmission<sup>9</sup>. Virulence in *V. cholerae* PG is  
41 mainly determined by two virulence factors: the cholera toxin (CT) and the toxin-coregulated  
42 pilus (TCP), which are encoded within horizontally acquired genetic elements, the CTX $\Phi$  phage  
43 and the Vibrio Pathogenicity Island-1 (VPI-1) respectively<sup>10-12</sup>. Both CTX $\Phi$  and VPI-1 are always  
44 found in the PG group, however, they are also encoded in some environmental populations of  
45 *V. cholerae*<sup>13-15</sup>. Furthermore, even though some non-PG strains can cause gastrointestinal  
46 infections, only strains from the PG clade have ever emerged as a source of pandemic  
47 cholera<sup>16</sup>.

48 To investigate the evolutionary origins of pandemic clones of *V. cholerae* and the  
49 potential for their reemergence, we analyzed 43 *V. cholerae* genomes sequenced from clinical  
50 and environmental samples (Methods; Supplementary Table 1). These genomes span the  
51 known genetic diversity of *V. cholerae* (Supplementary Notes), and were divided into a primary  
52 dataset of 22 genomes and a replication dataset of 22 genomes, with one reference genome in  
53 common (Supplementary Table 2). In the primary dataset, we chose 7 PGs, including both

54 classical strains, the source of the first six pandemics, and El Tor, to represent the genetic  
55 diversity of the pandemic group. We compared these with 15 non-clinical environmental  
56 genomes (EGs): 10 EGs from worldwide samples to include global diversity, and five sympatric  
57 isolates from the Great Bay Estuary (GBE) in New Hampshire, USA, a region with no recent  
58 history of cholera outbreaks<sup>17</sup>.

59 Consistent with the results of previous studies<sup>7,9,17,18</sup>, PGs form a distinct monophyletic  
60 group compared to EGs, based on the aligned core genome (Fig. 1a). Other than the PG group,  
61 there is little phylogenetic structure and the tree is star-like in both datasets (Supplementary Fig.  
62 1). Reticulations in the phylogenetic network indicate recombination or homoplasies (repeated  
63 mutations in independent lineages at the same locus), consistent with a large, genetically  
64 diverse and recombining *V. cholerae* population<sup>18,19</sup> (Supplementary Fig. 1). Given such a  
65 recombining population with mobile virulence factors, it remains puzzling why the ability to  
66 cause pandemic cholera is limited to the PG group.

67 To find unique features of PGs that could explain their pandemic success, we first  
68 searched for genes present in PGs but absent in EGs. Although there are a few genes and  
69 gene clusters which appear to be universally present in PGs, including CTX $\Phi$  and VPI-1, none  
70 of these are unique to PGs as they are also present in some EGs (Supplementary Table 3).  
71 This observation is consistent with previous work showing that these virulence genes are rapidly  
72 gained and lost in both EGs and PGs<sup>7,18</sup>. Therefore, there appear to be no gene families whose  
73 presence can easily explain the origin of pandemic cholera, nor are there strong boundaries to  
74 gene transfer between PGs and EGs.

75 Given the lack of PG-specific genes, we hypothesized that the environmental ancestor of  
76 the PGs had a particular genomic background containing alleles of core genes – which we term  
77 virulence adaptive polymorphisms (VAPs) – that served as “preadaptations” and enhanced its  
78 potential to give rise to pandemic disease. We started our search for VAPs by identifying SNPs

79 in the aligned core of the 22 primary dataset genomes with one allele uniquely present in all  
80 PGs and a different allele uniquely present in all EGs. We called this a “fixed” SNP pattern. We  
81 identified 819 such fixed SNPs distributed across the genome (Fig. 1b; Supplementary Fig. 3).  
82 Some fixed SNPs could have contributed to the evolution of the pandemic phenotype in PG,  
83 while others could be selectively neutral hitchhikers on the PG genomic background. Using the  
84 McDonald-Kreitman test (Methods), we found evidence for genome-wide positive selection  
85 during the divergence of PGs from EGs, due to an excess of nonsynonymous changes at fixed  
86 sites, in both datasets (Supplementary Table 4). However, no individual gene showed evidence  
87 for selection after correcting for multiple tests (Methods), rendering it difficult to identify  
88 candidate VAPs using fixed SNPs. Nonetheless, fixed SNPs constitute only a modest fraction of  
89 possible SNP patterns (Supplementary Table 2), and VAPs could also exist at other SNP sites  
90 that might shed light in the evolutionary past of the ancestral PG.

91 We then defined and searched for a “mixed” SNP pattern, where PGs encode one fixed  
92 allele and EGs encode a “mix” of two alleles: one that is unique to EGs and also one that is  
93 fixed in PGs (Fig. 1c inset). We reasoned that the existence of PG-like alleles segregating in  
94 contemporary environmental populations of *V. cholerae* could be informative about: 1) pathogen  
95 emergence, because having an allele fixed in PG suggests that it could also have been present  
96 in the environmental ancestor of PG; and 2) the potential for pathogen reemergence, as PG-like  
97 alleles, and thus potential VAPs, are still circulating in the environmental gene pool. We  
98 identified 39,171 mixed SNPs in the primary dataset. Most genes contain few mixed SNPs  
99 (median of 3) but some genes contain dense clusters, resulting in a mean of 10.3 mixed SNPs  
100 per gene. The replication dataset contained greater genetic diversity, but showed the same  
101 pattern of a few genes containing many mixed SNPs (Supplementary Table 2). Clusters of  
102 genes containing many mixed SNPs are visible when plotted across the genome (Fig. 1c).  
103 Some of these clusters are known polymorphic regions of the genome and could be mutation

104 hotspots containing SNPs not directly relevant to virulence adaptation. However, clusters of  
105 mixed SNPs do not visibly overlap with clusters of overall polymorphism (Fig. 1c), indicating that  
106 accumulation of mixed SNPs cannot be explained only by mutation hotspots.

107 To formally exclude mutation hotspots and focus on clusters of mixed SNPs shaped  
108 mainly by natural selection instead of mutation and drift, we considered only genes with an  
109 excess of nonsynonymous (NS) mixed SNPs relative to synonymous (S) mixed SNPs, which  
110 control for the baseline mutation rate. The mixed SNP pattern by definition groups PG-like EGs  
111 away from the other environmental strains and clusters them with the PGs. We reasoned that  
112 genes with an elevated NS:S ratio at mixed SNP sites were more likely to show phenotypic  
113 variations and have evolved under positive selection, possibly underlying preadaptations of the  
114 PG ancestor. Using a threshold of the number of mixed NS SNPs, NS:S, and  $dN/dS$  all two  
115 standard deviations above their genome-wide medians (Methods), we identified five genes as  
116 candidate VAPs in the primary dataset, three of which survived multiple hypothesis correction,  
117 and two of which (*ompU* and hypothetical gene VCD\_001600) were also found in the replication  
118 dataset (Table 1). In contrast to the star-like genome-wide phylogeny (Fig. 1a), each of these  
119 five gene trees support one or more EGs grouping with PGs (Fig. 1d and Supplementary Fig. 4).  
120 Three additional VAPs – all hypothetical proteins – were identified in the replication dataset,  
121 suggesting the potential for other candidate VAPs to be identified with further sampling of  
122 genetically diverse environmental genomes (Supplementary Table 5).

123 Among the candidate VAPs, the gene with the most significant excess of mixed  
124 nonsynonymous SNPs in both datasets is *ompU* (Table 1). OmpU is an outer membrane porin  
125 that has been shown to play numerous roles in the intestinal colonization of *V. cholerae*, making  
126 it a compelling candidate for phenotypic characterization<sup>20-23</sup>. We observed that environmental  
127 strains RC385, GBE0658 and GBE0428 have PG-like *ompU* alleles whereas most  
128 environmental strains, such as GBE1114, branch distantly from PG (Fig. 1d). We hypothesized

129 that the PG-like alleles present in environmental strains might confer properties conducive to  
130 virulence. To test this, we constructed three different mutant strains each encoding one of three  
131 environmental alleles of *ompU* (all from sympatric GBE strains) into the background of N16961,  
132 a clinical strain from the current pandemic (Fig. 2). OmpU was detected on a protein gel stained  
133 with Coomassie blue and through immunoblot in all the constructed strains, indicating that all  
134 the strains effectively produce the environmental versions of the protein (Fig. 2a and  
135 Supplementary Fig. 6). The band size differences reflect the diversity in protein sizes among the  
136 three environmental OmpU variants (Supplementary Fig. 5 and Supplementary Table 6).

137 We performed three sets of experiments to compare the phenotypes conferred by EG  
138 and PG-like alleles of *ompU*. First, we determined the survival of these strains in the presence  
139 of 0.4% bile, as it has been previously shown that OmpU confers resistance to this antimicrobial  
140 compound<sup>20</sup>. The mutant strain encoding the *ompU* allele from the environmental strain  
141 GBE1114 (OmpU<sup>GBE1114</sup>) showed a similar survival rate in the presence of bile as a deletion  
142 mutant (Fig 2b). In contrast, OmpU<sup>GBE0658</sup> show similar survival in the presence of bile as a PG  
143 strain (Fig. 2b). These experiments indicate that some environmental alleles of *ompU* confer  
144 properties beneficial for virulence. Second, we tested the survival of the mutants in the presence  
145 of polymyxin B, as OmpU also confers resistance against this antibiotic<sup>21</sup>. The ability to tolerate  
146 the antimicrobial effects of polymyxin B appears to be independent of which *ompU* allele is  
147 encoded by *V. cholerae*, as the three strains encoding environmental alleles of *ompU* had a  
148 similar survival rate (Fig. 2c). This experiment shows that OmpU<sup>GBE1114</sup> is not simply a loss of  
149 function mutant, and is not equivalent to a knockout. Third, we determined the intestinal  
150 colonization of the *ompU* mutants by performing competition assays using the infant mouse  
151 model of human infection. We found that OmpU<sup>GBE1114</sup> had a colonization defect similar to  
152  $\Delta ompU$  whereas OmpU<sup>GBE0658</sup> was able to colonize similarly to the strain with the wild-type PG

153 allele (Fig. 2d). These results indicate that certain naturally occurring environmental alleles of  
154 *ompU* confer properties that provide an advantage to *V. cholerae* during or prior to host  
155 colonization, as would be expected for VAPs.

156 The *ompU*<sup>GBE1114</sup> allele appears to be maladaptive for intestinal colonization; however, its  
157 presence in several environmental isolates of *V. cholerae* prompted us to investigate its  
158 possible adaptive role in the environment (Fig. 2e). *V. cholerae* forms biofilms on the surface of  
159 biotic and abiotic environmental surfaces<sup>24-26</sup>, yet biofilm formation inside the host is thought to  
160 impair intestinal colonization<sup>22,27</sup>. Strains with deletions in *ompU* have been shown to form a  
161 more robust biofilm on abiotic surfaces<sup>25</sup>. We found that OmpU<sup>GBE1114</sup> has higher biofilm  
162 formation than the strain encoding the wild-type PG allele, similar to the  $\Delta ompU$  strain (Fig. 2e).

163 Both OmpU<sup>GBE0658</sup> and OmpU<sup>GBE0428</sup> formed biofilm similar to the strains with the wild-type PG  
164 allele (Fig. 2e). It therefore appears that there is an evolutionary trade-off between encoding the  
165 PG-like or EG-like alleles of VAPs, as they seem to confer mutually exclusive traits: either  
166 biofilm formation or bile resistance and host colonization. This suggests that environmental  
167 strains can be divided into subgroups which, due to their contrasting lifestyles, differ in their  
168 potential to give rise to pandemic clones.

169 We have determined that virulence adaptive polymorphisms are present in the  
170 environment, and shown how these VAPs can be identified, based on two independent sets of  
171 genomes. The top candidate VAP, *ompU*, was identified in both of our genomic datasets. Our  
172 experiments show that the *ompU* allele from some environmental strains, such as GBE0658,  
173 confers properties that allow for host colonization equally as efficient as alleles from clinical  
174 strains (Fig. 2). This leads to a natural question: Why have environmental strains with PG-like  
175 alleles not emerged as pandemic cholera strains? It appears that a variety of virulence adaptive



176 genes and alleles are circulating in the environment, but only the PG group encodes the optimal  
177 combination of VAPs that allowed for pandemic potential (Fig. 3). We propose a conceptual  
178 model in which VAPs circulate in a diverse, recombining environmental gene pool, being  
179 maintained in the population through various biotic and abiotic selective pressures (Fig. 3a). A  
180 new ecological opportunity occurs, such as human consumption of brackish water or transient  
181 colonization of other animal hosts, which leads to the proliferation and gradual enrichment in the  
182 population of clones encoding a mosaic of VAPs (Fig. 3b). Finally, a genome encoding a critical  
183 combination of VAPs acquires key virulence factors allowing it to emerge as a virulent,  
184 potentially pandemic clone (Fig. 3c).

185 Our model posits that VAPs are circulating in the environment prior to the acquisition of  
186 key virulence factors. This is based on experimental evidence that a current pandemic strain  
187 cannot efficiently colonize the mammalian intestine without a PG-like *ompU* allele. If virulence  
188 adaptive alleles of *ompU* are indeed required prior to the acquisition of virulence factors, we  
189 would expect the same phenotypes of PG-like and EG-like *ompU* alleles in the genomic  
190 background of a more deeply branching PG isolate, such as classical *V. cholerae*. Indeed, we  
191 found that PG-like *ompU* alleles in the classical O395 background conferred efficient host  
192 colonization (Supplementary Fig. 7), which is consistent with an *ompU* VAP having played a role  
193 in the emergence of pandemic *V. cholerae*.

194 Our model further postulates that virulence adaptive alleles become enriched in the  
195 environmental population. Such enrichment would be made possible if these alleles provided a  
196 selective advantage in a newly available ecological niche, such as a human population  
197 consuming brackish water<sup>28</sup>. In previous work, we modeled an evolving, recombining microbial  
198 population that encounters a new ecological opportunity<sup>29</sup>. When adaptation to the new niche  
199 depends on few loci under positive selection in the niche, it is more likely for recombination to  
200 assemble the right combination of alleles in the same genome. These loci (akin to VAPs) could

201 contribute additively or synergistically to fitness. For example, the PG-like *ompU* alleles confer a  
202 tenfold increase in fitness during host colonization (Fig. 2d). Other VAPs might contribute further  
203 to this enhancement in host colonization. As more loci are involved in adaptation, it becomes  
204 less likely to achieve the optimal combination. Assuming the *V. cholerae* population undergoes  
205 approximately 100 recombination events per locus per generation (6.5 recombination events for  
206 every point mutation)<sup>19</sup>, equivalent to a recombination rate of  $10^{-4}$  in the modelled population of  
207 size  $10^6$ , an optimal combination of alleles at five loci could conceivably evolve, but seven loci is  
208 very unlikely<sup>29</sup>. Therefore, if virulence depended on five or fewer positively selected loci in the *V.*  
209 *cholerae* genome, the optimal combination of alleles would be expected to appear repeatedly in  
210 nature. Given suitable ecological opportunities, it is then plausible that pandemic *V. cholerae*  
211 could emerge multiple times, originating from outside the PG group. However, the number of  
212 loci that are sufficient for the emergence of a virulent strain remains unknown and if it was much  
213 greater than five, pathogen emergence would be naturally limited. We identified eight candidate  
214 VAPs in our two datasets, and if all eight are phenotypically confirmed to be VAPs, this number  
215 of VAPs (>5) might naturally limit pandemic clone emergence. Furthermore, these eight  
216 candidate VAPs passed stringent filters and we suspect there might exist additional VAPs in the  
217 genome, identifiable by further sampling and experimentation.

218         Here we have described a framework for identifying loci that are present in a natural  
219 population and confer properties beneficial for virulence prior to acquisition of essential  
220 virulence genes and host colonization. This framework could be applied to other bacterial  
221 pathogens that emerge as clonal offshoots from non-virulent relatives, including *Yersinia*,  
222 *Salmonella*, *Escherichia*, or other pathogenic *Vibrio* species<sup>1-4,6</sup>. Pathogens that emerged  
223 through clonal expansion limit our ability to dissect the genetic basis of their pathogenicity,  
224 because bacterial genome-wide association studies lack power when the phenotype of interest  
225 has evolved only once<sup>30</sup>. Our framework therefore provides a way forward to identify the genetic

- 226 basis of virulence, even in pathogens that evolved through clonal expansion, and begin to
- 227 assess the risk of pathogen emergence and reemergence from environmental gene pools.

228 **Methods**

229 **Genome sequencing.** DNA from clinical isolates (Bgd1, Bgd5, Bgd8, MQ1795<sup>31,32</sup>) and  
230 environmental isolates (GBE0428, GBE0658, GBE1068, GBE1114, GBE1173<sup>17</sup>) was extracted  
231 using the Gentra kit (QIAGEN) and purified using the MoBio PowerClean Pro DNA Clean-Up  
232 Kit. Multiplexed genomic libraries were constructed using the Illumina-compatible Nextera DNA  
233 Sample Prep kit following the manufacturer's instructions. Sequencing was performed with 250-  
234 bp paired-end (v2 kit) reads on the illumina MiSeq.

235

236 **Genome assembly.** To exclude low-quality data, raw reads were filtered with Trimmomatic<sup>33</sup>.  
237 The 15 first bases of each reads were trimmed and reads containing at least one base with a  
238 quality score of <30 were removed. *De novo* assembly was performed on the resulting reads  
239 using Ray v2.3.1<sup>34</sup>.

240

241 **Genome alignment, annotation and SNP calling.** We used mugsy v.1 r.2.2<sup>35</sup> with default  
242 parameters to align the primary dataset of 22 *V. cholerae* genomes (Supplementary Table 1).  
243 From this alignment, we extracted dimorphic SNP sites and annotated genes according to MJ-  
244 1236 as a reference genome. We defined the core genome as locally colinear blocks (LCBs)  
245 with all 22 genomes present in the alignment. We replicated the alignment, annotation and SNP  
246 calling using 21 different *V. cholerae* genomes, mainly from Orata *et al.*<sup>18</sup>, plus the MJ-126  
247 reference (Supplementary Table 1).

248

249 **Definition of orthologous groups.** Genomes were annotated using the RAST web server  
250 ([www.rast.nmpdr.org](http://www.rast.nmpdr.org))<sup>36</sup>. Annotated genes were clustered into orthologous groups using  
251 OrthoMCL ([www.orthomcl.org](http://www.orthomcl.org))<sup>37</sup> with default parameters, yielding 2844 orthologous groups.

252  
253 **Phylogenetic analysis.** We constructed a core genome phylogeny using the concatenated  
254 alignment of 1031 single-copy orthologous protein-coding genes (present in exactly one copy in  
255 each of the 22 primary dataset genomes). Each protein sequence was aligned with Muscle<sup>38</sup>,  
256 and the concatenated alignment was used to infer an approximate maximum likelihood  
257 phylogeny with FastTree v. 2.1.8<sup>39</sup> using default parameters (Fig. 1a). Individual gene trees (Fig.  
258 1d) were built in the same way. We constructed a neighbour-net of the 22 genomes using  
259 SplitsTree v.4.10<sup>40</sup>, based on dimorphic SNPs from the mugsy genome alignment, excluding  
260 sites with gaps.

261  
262 **Tests for selection.** We conducted a genome-wide version of the McDonald-Kreitman test<sup>41</sup> by  
263 first counting the number of fixed nonsynonymous (fn), fixed synonymous (fs), polymorphic  
264 nonsynonymous (pn), and polymorphic synonymous (ps) sites within each gene. We then  
265 summed these values across all genes (FN, FS, PN, and PS) and calculated the genome-wide  
266 Fixation Index,  $FI=(FN/FS)/(PN/PS)$ . A fixation index greater than one suggests positive  
267 selection between the ingroup and outgroup (in this case, between EGs and PGs). However,  
268 care must be taken when computing a genome-wide FI because summing genes with different  
269 amounts of substitutions and polymorphism can result in  $FI>1$  in the absence of selection<sup>42</sup>. We  
270 therefore performed 1000 permutations of the data, keeping the row totals (fn+fs and pn+ps)  
271 and column totals (fn+pn and fs+ps) constant and recomputing FI. We used the mean FI from  
272 the permutations as the expected value of FI under neutral evolution. To evaluate the  
273 hypothesis that the observed FI was higher than expected, suggesting positive selection, we  
274 computed a *P*-value as the fraction of permutations with FI greater or equal to the observed FI.

275 We repeated the test using polymorphism from either the PG group or the EG group  
276 (Supplementary Table 3).

277 To identify individual genes under selection between PGs and EGs in the primary  
278 dataset, we restricted our search to 87 genes with  $fn > 0.68$  and  $fn:fs > 1.48$  (respectively two  
279 standard deviations about the genome-wide medians). We then calculated the gene-specific FI  
280 and assessed its significance with a Fisher exact test. We found no genes with FI significantly  
281 greater than one, after Bonferroni correction for 87 tests. Similarly, in the replication dataset, we  
282 restricted our search to 26 genes with  $fn > 1.5$  and  $fn:fs > 1.70$ , none of which had significantly  
283 high FI after correction for multiple tests.

284 To identify genes with an excess of nonsynonymous mixed SNPs (likely due to selection  
285 for amino acid changes), we restricted our search to five genes with  $\geq 12$  NS mixed SNPs per  
286 gene and mixed NS:S  $> 1.78$  (respectively two standard deviations above the genome-wide  
287 medians). We used a one-sided binomial test to assess whether the observed NS:S ratio for  
288 each gene was significantly greater than the genome-wide median NS:S ratio of 0.5 (after  
289 adding a pseudocount of one to both NS and S). Three out of the five genes had a significantly  
290 high mixed NS:S ratio ( $P < 0.05$ ) after Bonferroni correction for five tests (Table 1). We repeated  
291 this procedure in the replication dataset, identifying genes with  $\geq 18$  NS mixed SNPs per gene  
292 and mixed NS:S  $> 1.65$  (respectively two standard deviations above the genome-wide  
293 medians). We used a one-sided binomial test to assess whether the observed NS:S ratio for

294 each gene was significantly greater than the genome-wide median NS:S ratio of 0.33 (after  
295 adding a pseudocount of one to both NS and S). The results of these tests are shown for genes  
296 also identified in the primary dataset (Table 1) and three additional genes identified in the  
297 replication dataset (Supplementary Table 5).

298

299 **Bacterial strains and plasmids.** *V. cholerae* O395 and *V. cholerae* N16961 were used as wild-  
300 type strains of classical and El Tor biotypes respectively. Strains cultivated on solid medium  
301 were grown on LB agar; strains in liquid media were grown in aerated LB broth at 37°C.  
302 pKAS154 was used for allelic exchange<sup>43</sup>. When necessary, media was supplemented with  
303 antibiotics to select for certain plasmids or strains of *V. cholerae* at the following concentrations:  
304 gentamycin, 30µg/ml; kanamycin, 45µg/ml; polymyxin B, 50µg/ml; and streptomycin, 1 mg/ml.

305

306 **Strain construction.** In-frame deletions and exchange of *ompU* alleles in both O395 and  
307 N16961 biotypes were constructed via homologous recombination<sup>43</sup>. For *ompU* deletions, PCR  
308 was used to amplify two 500 bp fragments flanking the *ompU* gene and to introduce restriction  
309 sites for cloning into the suicide vector pKAS154. For exchange of environmental *ompU* alleles,  
310 the respective allele was also amplified using the extracted DNA from each environmental  
311 strain. Restriction sites were introduced in the primers. The fragments were then cloned into a  
312 restriction-digested suicide plasmid, pKAS154, using a four-segment ligation for each  
313 environmental allele exchange mutants (plasmid, *ompU* flanking fragments and environmental  
314 *ompU* allele). The resulting plasmids were electroporated into *Escherichia coli* S17-1λpir. *E. coli*  
315 with the constructed inserts. Different sizes of *ompU* genes were confirmed, sequenced and

316 compared to the genome sequences of environmental strains. For both deletion and exchange  
317 mutants, plasmids with the insert of interest were mated with wild-type *V. cholerae* O395 or  
318 N16961, and allelic exchange was carried out by selection on antibiotics as described  
319 previously<sup>43</sup>. For a more detailed description of allelic exchange please refer to Skorupski and  
320 Taylor 1998<sup>43</sup>. Potential mutants were screened using PCR: two primers flanking the deletion  
321 construct were used to amplify chromosomal DNA isolated from plated *V. cholerae*. The lengths  
322 of the PCR fragments were analyzed on 0.8% agarose gel for gene deletions and putative  
323 deletions and allele exchange were subsequently confirmed by DNA sequencing.

324

325 **OmpU visualization and immunoblotting.** Whole cell protein extracts were prepared from  
326 cultures grown for overnight at 37°C in a rotary shaker. The extracts were subjected to SDS-  
327 PAGE on 16% Tris Glycine gels (Invitrogen). OmpU bands were visualized after protein gels  
328 were stained by Coomassie blue overnight. Prior to staining the gels were transferred to  
329 nitrocellulose membranes using iBlot (Invitrogen). The membranes were blocked O/N in Tris-  
330 Buffered Saline, 3% BSA. Primary OmpU antibodies were diluted 1:10,000 in TBST (Tris-  
331 Buffered Saline, 0.5% Tween-20). Membranes were incubated with primary antibodies for 2  
332 hours at room temperature. After incubation, the membranes were washed with TBST four  
333 times. Goat anti-rabbit secondary antibodies (BioRad) were diluted 1:10,000 in TBST and  
334 incubated for 30 minutes at room temperature. The membranes were washed 4 times with TBS  
335 (Tris-Buffered Saline). Reactive protein bands were detected via ECL (Amersham).

336

337 **Survival assays.** *V. cholerae* strains were cultured overnight in LB broth at 37°C in a rotary  
338 shaker. Overnight cultures were diluted 1:100 in LB and grown to an OD600 of 0.5. Cells were



339 pelleted and resuspended in either LB broth, LB containing 0.4% bile bovine (Sigma), or LB  
340 containing 1000U/ml of polymyxin B (Sigma). Cultures were incubated for 1h at 37°C in a rotary  
341 shaker. After incubation CFU/ml of each culture was calculated by plating dilutions in LB plates.  
342 Survival was calculated by comparing the number of CFU/ml in LB plus treatment versus LB.  
343  $N \geq 6$ . No samples were excluded.

344

345 **Infant mouse competition assays.** Overnight cultures were diluted 1:100. Each test strain was  
346 mixed in a 1:1 ratio with a  $\Delta lacZ$  reference strain. Four- to five-day-old CD-1 mice (*Mus*  
347 *musculus*) from several mixed litters were randomly inoculated orogastrically in blinded  
348 experiments with 50 $\mu$ l of the bacterial mixture. Sex of the animals was not inspected prior to  
349 inoculations. The intestines were harvested 24h post-inoculation and homogenized in 4ml of LB  
350 broth containing 10% glycerol. The mixtures were serially diluted and plated on LB agar plates  
351 supplemented with streptomycin and 5-bromo-4-chloro-3-indolyl-D-galactopyramoside (X-Gal)  
352 (40 $\mu$ g/ml). The competition indices were calculated as previously described by others, test  
353 (output CFUs/input CFUs)/reference (output CFUs/ input CFUs) and the sample size selected  
354 was appropriate for statistical analysis. No samples or animals were excluded. Animal work was  
355 approved by the Institutional Animal Care and Use Committee (IACUC).

356

357 **Biofilm assays.** 96-well plate assay. Cultures were incubated overnight at 30°C. 100 $\mu$ l of 1:100  
358 dilutions of overnight cultures were placed per well in 96-well plates. Plates were left at 25°C for  
359 24h. Liquid contents were discarded and plates were washed 2 times with LB. 200 $\mu$ l of 0.01%

360 crystal violet was added per well and incubated at room temperature for 5 minutes. Liquid  
361 contents were discarded and plates were washed extensively with dH<sub>2</sub>O. After the plates were  
362 dry, biofilms were resuspended in 150µl of 50% acetic acid. Contents were transferred to a flat  
363 bottom dish and quantitated in a microtiter plate reader at OD550. Values were plotted using  
364 Prism software. N=15.

365 **Data availability**

366 The nine genomes sequenced in this study (Bgd1, Bgd5, Bgd8, MQ1795, GBE0428, GBE0658,  
367 GBE1068, GBE1114, GBE1173) have been deposited under BioProject ID PRJNA349157 in  
368 NCBI GenBank under accession numbers SAMN05924900-SAMN05924908 (Supplementary  
369 Table 1). All other data supporting the findings of this study are available from the authors upon  
370 request.

371

## 372 References

- 373 1. McNally, A., Thomson, N. R., Reuter, S. & Wren, B. W. 'Add, stir and reduce': *Yersinia* spp. as model  
374 bacteria for pathogen evolution. *Nat. Rev. Microbiol.* **14**, 177–190 (2016).
- 375 2. Luo, C. *et al.* Genome sequencing of environmental *Escherichia coli* expands understanding of the ecology  
376 and speciation of the model bacterial species. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 7200–7205 (2011).
- 377 3. Reen, F. J., Almagro-Moreno, S., Ussery, D. & Boyd, E. F. The genomic code: inferring Vibrionaceae niche  
378 specialization. *Nat. Rev. Microbiol.* **4**, 697–704 (2006).
- 379 4. Ochman, H. Genes Lost and Genes Found: Evolution of Bacterial Pathogenesis and Symbiosis. *Science* **292**,  
380 1096–1099 (2001).
- 381 5. Shapiro, B. J. How clonal are bacteria over time? *Curr. Opin. Microbiol.* **31**, 116–123 (2016).
- 382 6. Cui, Y. *et al.* Epidemic Clones, Oceanic Gene Pools, and Eco-LD in the Free Living Marine Pathogen *Vibrio*  
383 *parahaemolyticus*. *Mol. Biol. Evol.* **32**, 1396–1410 (2015).
- 384 7. Chun, J. *et al.* Comparative genomics reveals mechanism for short-term and long-term clonal transitions in  
385 pandemic *Vibrio cholerae*. *Proc. Natl. Acad. Sci. U.S.A.* **106**, 15442–15447 (2009).
- 386 8. Boucher, Y. Sustained local diversity of *Vibrio cholerae* O1 biotypes in a previously cholera-free country.  
387 *mBio* **7**, e00570–16 (2016).
- 388 9. Mutreja, A. *et al.* Evidence for several waves of global transmission in the seventh cholera pandemic. *Nature*  
389 **477**, 462–465 (2011).
- 390 10. Taylor, R. K., Miller, V. L., Furlong, D. B. & Mekalanos, J. J. Use of *phoA* gene fusions to identify a pilus  
391 colonization factor coordinately regulated with cholera toxin. *Proc. Natl. Acad. Sci. U.S.A.* **84**, 2833–2837  
392 (1987).
- 393 11. Waldor, M. K. & Mekalanos, J. J. Lysogenic conversion by a filamentous phage encoding cholera toxin.  
394 *Science (New York, N.Y.)* **272**, 1910–1914 (1996).
- 395 12. Karaolis, D. K. *et al.* A *Vibrio cholerae* pathogenicity island associated with epidemic and pandemic strains.  
396 *Proc. Natl. Acad. Sci. U.S.A.* **95**, 3134–3139 (1998).
- 397 13. Faruque, S. M. *et al.* Genetic diversity and virulence potential of environmental *Vibrio cholerae* population in a  
398 cholera-endemic area. *Proc. Natl. Acad. Sci. U.S.A.* **101**, 2123–2128 (2004).
- 399 14. Rivera, I. N., Chun, J., Huq, A., Sack, R. B. & Colwell, R. R. Genotypes associated with virulence in  
400 environmental isolates of *Vibrio cholerae*. *Appl. Environ. Microbiol.* **67**, 2421–2429 (2001).
- 401 15. Gennari, M., Ghidini, V. & Lleo, M. M. Virulence genes and pathogenicity islands in environmental *Vibrio*  
402 strains non-pathogenic to humans. *FEMS Microbiol. Ecol.* **82**, 563–573 (2012).
- 403 16. Dziejman, M. *et al.* Genomic characterization of non-O1, non-O139 *Vibrio cholerae* reveals genes for a type  
404 III secretion system. *Proc. Natl. Acad. Sci. U.S.A.* **102**, 3465–3470 (2005).
- 405 17. Schuster, B. M. *et al.* Ecology and genetic structure of a northern temperate *Vibrio cholerae* population  
406 related to toxigenic isolates. *Appl. Environ. Microbiol.* **77**, 7568–7575 (2011).
- 407 18. Orata, F. D. *et al.* The dynamics of genetic interactions between *Vibrio metoecus* and *Vibrio cholerae*, two  
408 close relatives co-occurring in the environment. *Genome Biol Evol* **7**, 2941–2954 (2015).
- 409 19. Keymer, D. P. & Boehm, A. B. Recombination shapes the structure of an environmental *Vibrio cholerae*  
410 population. *Appl. Environ. Microbiol.* **77**, 537–544 (2011).
- 411 20. Provenzano, D., Schuhmacher, D. A., Barker, J. L. & Klose, K. E. The virulence regulatory protein ToxR  
412 mediates enhanced bile resistance in *Vibrio cholerae* and other pathogenic *Vibrio* species. *Infect. Immun.* **68**,  
413 1491–1497 (2000).
- 414 21. Mathur, J. & Waldor, M. K. The *Vibrio cholerae* ToxR-regulated porin OmpU confers resistance to  
415 antimicrobial peptides. *Infect. Immun.* **72**, 3577–3583 (2004).
- 416 22. Almagro-Moreno, S., Pruss, K. & Taylor, R. K. Intestinal colonization dynamics of *Vibrio cholerae*. *PLoS*  
417 *Pathog.* **11**, e1004787 (2015).
- 418 23. Merrell, D. S., Bailey, C., Kaper, J. B. & Camilli, A. The ToxR-mediated organic acid tolerance response of  
419 *Vibrio cholerae* requires OmpU. *J. Bacteriol.* **183**, 2746–2754 (2001).
- 420 24. Watnick, P. I. & Kolter, R. Steps in the development of a *Vibrio cholerae* El Tor biofilm. *Mol. Microbiol.* **34**,  
421 586–595 (1999).
- 422 25. Valeru, S. P., Wai, S. N., Saeed, A., Sandström, G. & Abd, H. ToxR of *Vibrio cholerae* affects biofilm, rugosity  
423 and survival with *Acanthamoeba castellanii*. *BMC Res Notes* **5**, 33 (2012).
- 424 26. Yildiz, F. H. & Schoolnik, G. K. *Vibrio cholerae* O1 El Tor: identification of a gene cluster required for the  
425 rugose colony type, exopolysaccharide production, chlorine resistance, and biofilm formation. *Proc. Natl.*  
426 *Acad. Sci. U.S.A.* **96**, 4028–4033 (1999).

- 427 27. Hsiao, A., Liu, Z., Joelsson, A. & Zhu, J. *Vibrio cholerae* virulence regulator-coordinated evasion of host  
428 immunity. *Proc. Natl. Acad. Sci. U.S.A.* **103**, 14542–14547 (2006).
- 429 28. Boucher, Y., Orata, F. D. & Alam, M. The out-of-the-delta hypothesis: dense human populations in low-lying  
430 river deltas served as agents for the evolution of a deadly pathogen. *Front. Microbiol.* **6**, L19401 (2015).
- 431 29. Friedman, J., Alm, E. J. & Shapiro, B. J. Sympatric Speciation: When Is It Possible in Bacteria? *PLoS ONE* **8**,  
432 e53539 (2013).
- 433 30. Chen, P. E. & Shapiro, B. J. The advent of genome-wide association studies for bacteria. *Curr. Opin.*  
434 *Microbiol.* **25**, 17–24 (2015).
- 435 31. Nair, G. B. *et al.* New variants of *Vibrio cholerae* O1 biotype El Tor with attributes of the classical biotype from  
436 hospitalized patients with acute diarrhea in Bangladesh. *J. Clin. Microbiol.* **40**, 3296–3299 (2002).
- 437 32. Son, M. S., Megli, C. J., Kovacicova, G., Qadri, F. & Taylor, R. K. Characterization of *Vibrio cholerae* O1 El  
438 Tor biotype variant clinical isolates from Bangladesh and Haiti, including a molecular genetic analysis of  
439 virulence genes. *J. Clin. Microbiol.* **49**, 3739–3749 (2011).
- 440 33. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data.  
441 *Bioinformatics* **30**, 2114–2120 (2014).
- 442 34. Boisvert, S., Raymond, F., Godzaridis, É., Laviolette, F. & Corbeil, J. Ray Meta: scalable de novo  
443 metagenome assembly and profiling. *Genome Biol* **13**, R122 (2012).
- 444 35. Samuel V Angiuoli, S. L. S. Mugsy: fast multiple alignment of closely related whole genomes. *Bioinformatics*  
445 **27**, 334–342 (2011).
- 446 36. Aziz, R. K. *et al.* The RAST Server: Rapid Annotations using Subsystems Technology. *BMC Genomics* **9**, 1  
447 (2008).
- 448 37. Li, L., Stoeckert, C. J. & Roos, D. S. OrthoMCL: identification of ortholog groups for eukaryotic genomes.  
449 *Genome Res.* **13**, 2178–2189 (2003).
- 450 38. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids*  
451 *Res.* **32**, 1792–1797 (2004).
- 452 39. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – Approximately Maximum-Likelihood Trees for Large  
453 Alignments. *PLoS ONE* **5**, e9490 (2010).
- 454 40. Bryant, D. & Moulton, V. Neighbor-net: an agglomerative method for the construction of phylogenetic  
455 networks. *Mol. Biol. Evol.* **21**, 255–265 (2004).
- 456 41. McDonald, J. H. & Kreitman, M. Adaptive protein evolution at the Adh locus in Drosophila. *Nature* **351**, 652–  
457 654 (1991).
- 458 42. Shapiro, J. A. *et al.* Adaptive genic evolution in the Drosophila genomes. *Proc. Natl. Acad. Sci. U.S.A.* **104**,  
459 2271–2276 (2007).
- 460 43. Skorupski, K. & Taylor, R. K. Positive selection vectors for allelic exchange. *Gene* **169**, 47–52 (1996).
- 461

462 **Correspondence should be addressed to:**

463 Salvador Almagro-Moreno (samoreno@ucf.edu)

464 **Acknowledgements**

465 The authors would like to thank the anonymous reviewers for their thoughtful comments and  
466 suggestions. We would also like to thank Otto Cordero, Yves Terrat, Nicolas Tromas and  
467 Britney Privett for constructive comments on the manuscript. We thank Lawrence Shelven for  
468 his highly valuable technical assistance. BJS was supported by a Canada Research Chair and  
469 the Canadian Institutes for Health Research. RKT was supported by a National Institutes of  
470 Health grants AI039654 and AI025096. SAM was supported by startup funds from the Burnett  
471 School of Biomedical Sciences at the University of Central Florida and Dartmouth College's E.  
472 E. Just Postdoctoral Fellowship. This article is dedicated to the memory of Ronald K. Taylor.

473 **Author contributions**

474 SAM conceived the study. BJS, RKT and SAM designed the study. IL sequenced genomes.

475 BJS performed computational analysis. GK and SAM performed phenotypic characterization.

476 BJS and SAM analyzed, interpreted data and wrote the article. All authors have read a

477 version of the manuscript.



478 **Figure Legends**

479 **Figure 1. Comparative genomics reveals candidate virulence adaptive polymorphisms. a,**

480 Phylogeny of 22 *V. cholerae* genomes based on 1031 single-copy orthologs in the primary  
481 dataset. All branches have local support values >0.99 (based on FastTree's approximate  
482 likelihood ratio test) except for very short, deep internal branches (resulting in the star-like  
483 polytomy at the centre of the tree). Not all 22 genomes are visible because some have nearly  
484 identical sequences (e.g. 6 of the 7 PG genomes are nearly identical, shown as an orange  
485 triangle; GBE1173 and GBE1114 are nearly identical, as can be seen in Supplementary Fig. 1).

486 **b,** Distribution of fixed SNPs across chromosome 1. (See Supplementary Fig. 3 for  
487 chromosome 2). Genome position is according to the MJ-1236 reference genome. SNP-free  
488 regions (e.g. near 3 Mbp, the locus of the integrative conjugative element) are part of the flexible  
489 genome, present in the reference but not the other 21 genomes. The schematic tree in the top  
490 left illustrates the fixed SNP pattern, in which one allele is present in PGs and a different allele  
491 in EGs. **c,** Distribution of mixed SNPs across the genome. The cartoon tree in the top left  
492 illustrates the mixed SNP pattern, in which one allele is fixed in PGs, and another allele is  
493 polymorphic among EGs, with some EGs containing the PG-like allele. Black arrows show  
494 candidate VAPs (Table 1). Grey arrow shows the flagellum as an example variable region not  
495 containing candidate VAPs. **d,** *ompU* phylogeny. All visible branches have local support values  
496 >0.9 except for the branch separating RC385 and GBE0658, the branch grouping MJ-1236 and  
497 O395 together, and the branch grouping HE09 and VL426 together.

498

499 **Figure 2. Phenotypic characterization of *ompU* alleles. a,** OmpU production in clinical strains  
500 of *V. cholerae* encoding environmental alleles of *ompU*. Total protein lysates were run on a 16%  
501 Tris-glycine gel. OmpU bands were visualized after protein gels were stained with Coomassie

502 blue. Three independent sets of protein lysates were examined and showed an identical band  
503 pattern. **b**, Survival of *ompU* mutants in the presence of bile (n=7) or **c**, polymyxin B (n=6). **d**,  
504 Colonization of the small intestine of *ompU* mutant strains (n=6). **e**, Biofilm formation of *ompU*  
505 mutant strains on an abiotic surface (n=15). Yellow bars and symbols, PG-like allele; red bars  
506 and squares,  $\Delta ompU$ ; blue bars and triangles, EG-like allele. Center values represent the mean  
507 and error bars the standard deviation. Variance between the groups was similar. Statistical  
508 comparisons were made using Student's *t*-test. \* $P < 0.05$ , \*\*\* $P < 0.001$

509

510 **Figure 3. Model of pandemic clone emergence from an environmental gene pool.** We  
511 propose a model that involves three events required for the emergence of pathogenic clones  
512 from environmental populations. **a**, selection of VAPs. Virulence adaptive alleles circulate in  
513 naturally occurring populations (orange symbols) and can be exchanged and mobilized through  
514 recombination (green dashed arrows). Ecological variation (e.g. temperature, nutrient  
515 availability, pH, etc.) leads to the selection of VAPs and an increase in their distribution in  
516 environmental populations. **b**, enrichment of clones. A new ecological opportunity occurs  
517 (human consumption of untreated waters, transient colonization of new environmental hosts,  
518 etc.) which leads to the proliferation and enrichment in the population of clones encoding a  
519 mosaic of VAPs. **c**, acquisition of virulence factors. A strain encoding a minimum set of VAPs  
520 required for host colonization acquires the virulence factors that are necessary to produce a  
521 successful infection and give rise to pandemic disease.

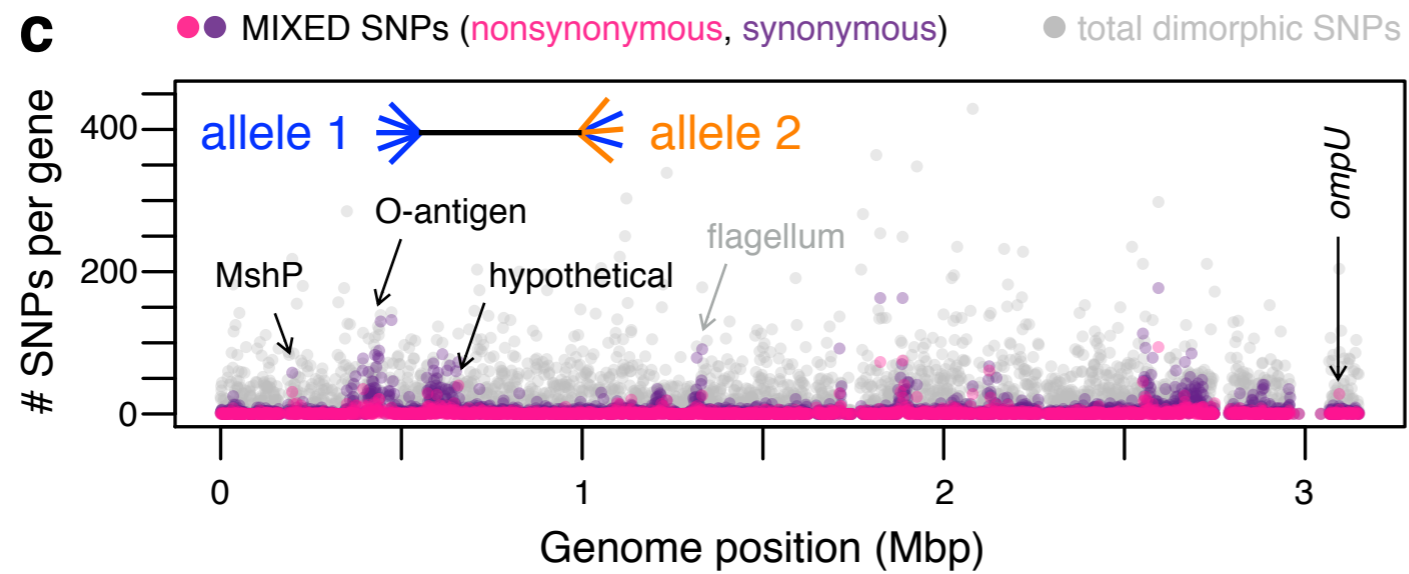
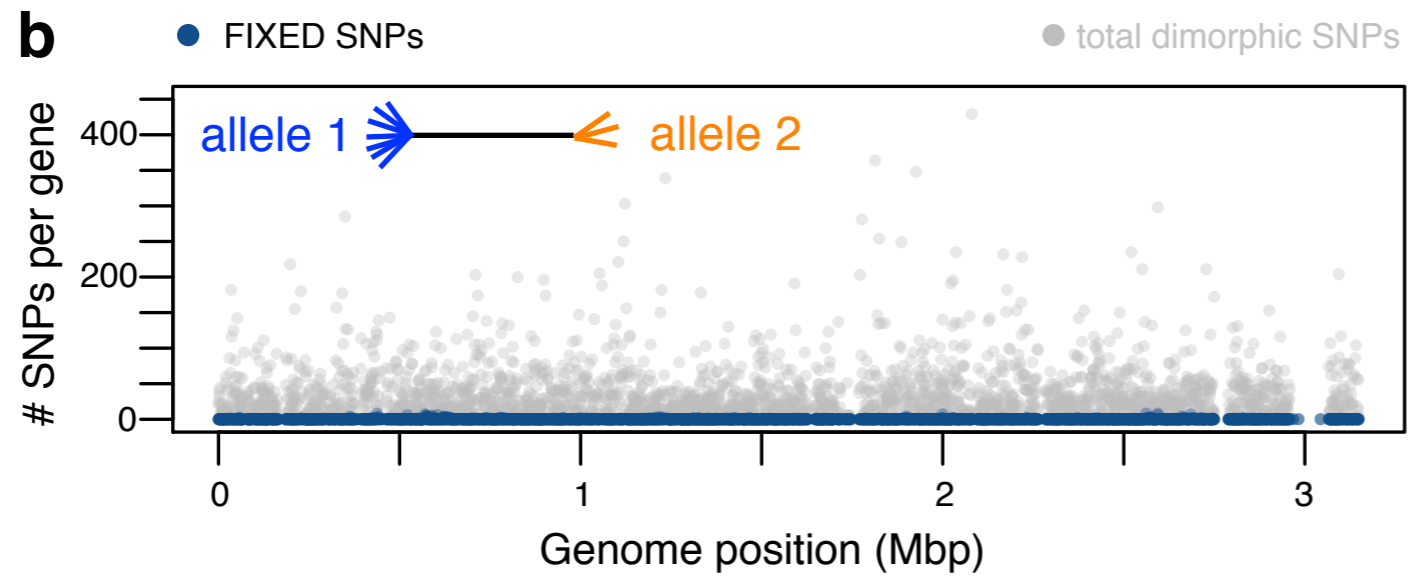
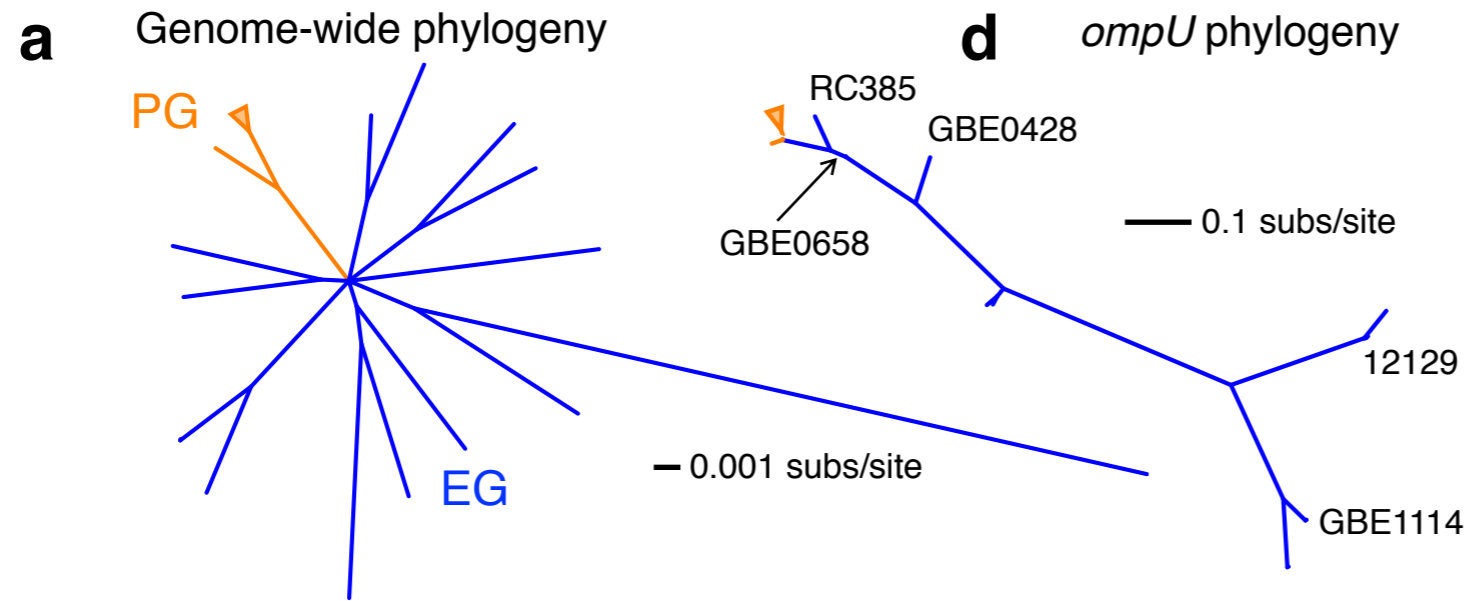
522

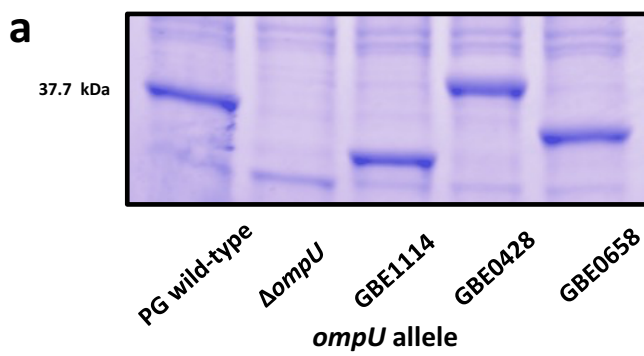
523 **Tables**

524 **Table 1. Characteristics of five predicted VAPs with an excess of nonsynonymous mixed**  
 525 **SNPs.**

Gene ID (VCD #)	Annotation	Gene length (bp)	Total # SNPs	Mixed <i>dN</i>	Mixed <i>dS</i>	Mixed <i>dN/dS</i>	Mixed NS	Mixed S	<i>P</i>
1 003778	outer membrane protein OmpU	1053	204 (110)	0.037 (0.018)	0.045 (0.033)	0.83 (0.53)	28 (13)	11 (8)	0.0047* (0.0068*)
2 001600	hypothetical	258	21 (26)	0.082 (0.12)	0.081 (0.065)	1.01 (1.82)	15 (22)	4 (3)	0.0096* (2.32e-8*)
3 001013	hypothetical	642	19 (13)	0.034 (0.006)	0.029 (0.012)	1.18 (0.55)	15 (2)	4 (1)	0.0096* (n.s.)
4 001209	MSHA biogenesis protein MshP	432	85 (63)	0.047 (0.025)	0.045 (0.054)	1.04 (0.46)	14 (7)	4 (5)	0.0154 (0.066)
5 001230	lipid A core O-antigen ligase	1794	75 (47)	0.0098 (0.001)	0.013 (0.004)	0.76 (0.18)	12 (0)	5 (1)	0.0717 (n.s.)

526 NS=nonsynonymous; S=synonymous. *dN* = number of nonsynonymous SNPs per  
 527 nonsynonymous site; *dS* = number of synonymous SNPs per synonymous site. Mixed SNPs are  
 528 polymorphic in EGs but fixed in PGs, meaning that at least one EG contains a PG-like allele.  
 529 The genes listed have mixed NS, mixed NS:S, and mixed *dN/dS* each over two standard  
 530 deviations above their respective genome-wide medians in the primary dataset. A one-sided  
 531 binomial test determined if the mixed NS:S ratio was greater than the expected genome-wide  
 532 median value of 0.5 per gene (uncorrected *P*-values shown; asterisks (\*) indicate *P* < 0.05 after  
 533 Bonferroni correction for five tests). Numbers in parentheses are for the replication dataset, with  
 534 an expected genome-wide median mixed N:S of 0.33. *P*-values greater than 0.1 are denoted as  
 535 not significant (n.s.). Genes 1, 2, 4, and 5 are indicated with arrows on Figure 1C. Gene 3 is on  
 536 chromosome 2 (Supplementary Fig. 3).





bioRxiv preprint doi: <https://doi.org/10.1101/063115>; this version posted December 13, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

