

1 Parallel evolution of two clades of a major Atlantic endemic *Vibrio parahaemolyticus* pathogen  
2 lineage by independent acquisition of related pathogenicity islands

3

4 Feng Xu<sup>1,2,3</sup>, Narjol Gonzalez-Escalona<sup>4</sup>, Kevin P. Drees<sup>1,2</sup>, Robert P. Sebra<sup>5</sup>, Vaughn S.  
5 Cooper<sup>1,2,\*</sup>, Stephen H. Jones<sup>1,6</sup>, and Cheryl A. Whistler<sup>1,2#</sup>

6

7 Running Title: parallel evolution of ST631 *Vibrio parahaemolyticus*

8

9 <sup>1</sup>Northeast Center for Vibrio Disease and Ecology, University of New Hampshire, Durham, NH;

10 <sup>2</sup>Department of Molecular, Cellular and Biomedical Sciences, University of New Hampshire,

11 Durham, NH; <sup>3</sup>Genetics Graduate Program, University of New Hampshire, Durham, NH;

12 <sup>4</sup>Center for Food Safety and Applied Nutrition, Food and Drug Administration, College Park,

13 MD; <sup>5</sup>Icahn Institute and Department of Genetics & Genomic Sciences, Icahn School of

14 Medicine at Mount Sinai, New York, NY; and <sup>6</sup>Department of Natural Resources and the

15 Environment, University of New Hampshire, Durham, NH, USA.

16

17 \*Current address: Microbiology and Molecular Genetics, University of Pittsburgh School

18 of Medicine, Pittsburgh, PA

19

20 #Corresponding author e-mail: [cheryl.whistler@unh.edu](mailto:cheryl.whistler@unh.edu)

## 21 **ABSTRACT**

22 Shellfish-transmitted *Vibrio parahaemolyticus* infections have recently increased from  
23 locations with historically low disease incidence, such as the Northeast United States (US). This  
24 change coincided with a bacterial population shift towards human pathogenic variants occurring  
25 in part through the introduction of several Pacific native strains (ST36, ST43 and ST636) to near-  
26 shore areas off the Atlantic coast of the Northeast US. Concomitantly, ST631 emerged as a  
27 major endemic pathogen. Phylogenetic trees of clinical and environmental isolates indicated that  
28 two clades diverged from a common ST631 ancestor, and in each of these clades, human  
29 pathogenic variants evolved independently through acquisition of distinct *Vibrio* pathogenicity  
30 islands (VPaI). These VPaI differ from each other and that bear little resemblance to hemolysin-  
31 containing VPaI from pandemic strains. Clade I ST631 isolates either harbored no hemolysins,  
32 or contained a chromosome I-inserted island we call VPaI $\beta$  that encodes a type three secretion  
33 system (T3SS2 $\beta$ ) typical of Trh hemolysin-producers. The more clinically prevalent and clonal  
34 ST631 clade II had an island we call VPaI $\gamma$  that encodes both *tdh* and *trh* and that was inserted in  
35 chromosome II. VPaI $\gamma$  was derived from VPaI $\beta$  but with some elements in common with  
36 pandemic strains, exemplifying the mosaic nature of pathogenicity islands. Genomics  
37 comparisons and amplicon assays identified VPaI $\gamma$ -type islands containing *tdh* inserted adjacent  
38 to the *ure* cluster in the three introduced Pacific and most other lineages that collectively cause  
39 67% of Northeast US infections as of 2016.

40

## 41 **IMPORTANCE**

42 The availability of three different hemolysin genotypes in the ST631 lineage provided a  
43 unique opportunity to employ genome comparisons to further our understanding of the processes

44 underlying pathogen evolution. The fact that two different pathogenic clades arose in parallel  
45 from the same potentially benign lineage by independent VP<sub>AI</sub> acquisition is surprising  
46 considering the historically low prevalence of community members harboring VP<sub>AI</sub> in waters  
47 along the Northeast US Coast that could serve as the source of this material. This illustrates a  
48 possible predisposition of some lineages to not only acquire foreign DNA but also to become  
49 human pathogens. Whereas the underlying cause for the expansion of *V. parahaemolyticus*  
50 strains harboring VP<sub>AI</sub> along the US Atlantic coast and spread of this element to multiple  
51 lineages that underlies disease emergence is not known, this work underscores the need to define  
52 the environment factors that favor strains harboring VP<sub>AI</sub> in locations of emergent disease.

53

## 54 INTRODUCTION

55 *Vibrio parahaemolyticus* is an emergent pathogen capable of causing human gastric  
56 infections when consumed, most often with contaminated shellfish (1, 2). Some human  
57 pathogenic *V. parahaemolyticus* strains evolve from diverse non-pathogenic communities  
58 through horizontal acquisition of *Vibrio* pathogenicity islands (VP<sub>AI</sub>) (3-5). Gastric pathogenic *V.*  
59 *parahaemolyticus* typically harbor islands with at least one of two types of horizontally acquired  
60 hemolysin genes (*tdh* and *trh*) that are routinely used for pathogen discrimination even though  
61 their role in disease appears modest (6-11). Most pathogenic *V. parahaemolyticus* isolates also  
62 produce accessory type three secretion systems (T3SS) that translocate effector proteins that  
63 contribute to host interaction (12-14). Two evolutionarily divergent horizontally-acquired  
64 accessory systems (T3SS2 $\alpha$  or T3SS2 $\beta$ ) contribute to human disease and are genetically linked  
65 to hemolysin genes (two *tdh* genes with T3SS2 $\alpha$ , and *trh* with T3SS2 $\beta$ ) in contiguous but distinct  
66 islands (4, 15-17). The first described *tdh*-harboring island [called by several different names

67 including Vp-PAI (15), VPai-7 (4), and *tdh*VPA (17)] from an Asian pandemic strain called  
68 RIMD 2210366 is fairly well-characterized (4, 5, 13, 18, 19). In contrast, islands containing  
69 T3SS2 $\beta$  linked to *trh* and a urease (*ure*) cluster, which confers a useful diagnostic phenotype,  
70 [where similar islands are described by others as Vp-PAI<sub>TH3966</sub> (16), or *trh*VPA(17, 20)] have  
71 received only modest attention. Strains harboring both *tdh* and *trh* are increasingly associated  
72 with disease in North America (21-26), and yet, to our knowledge, the exact configuration of  
73 hemolysin-associated VPai(s) in isolates that contain both *tdh* and *trh* have not yet been  
74 described [although see (20)]. Thus it is unclear how virulence loci and islands in these emergent  
75 pathogen lineages carrying both hemolysins evolved and spread.

76         The expanding populations of *V. parahaemolyticus* have increased infections even in  
77 temperate regions previously only rarely impacted by this pathogen and where most  
78 environmental isolates harbor no known virulence determinants (27). A related complex of Asia-  
79 derived pandemic strains, most often identified as serotype O3:K6 and also known as sequence  
80 type (ST) 3 (based on allele combinations of seven housekeeping genes) causes the most disease  
81 globally (28). An unrelated Pacific native strain called ST36 (also described as serotype O4:K12)  
82 currently dominates infections in North America, including from the Northeast United States  
83 (US) (21, 26, 29). The introduction of ST36 into the Atlantic Ocean by an unknown route  
84 precipitated a series of outbreaks from Atlantic shellfish starting in 2012 (29, 30). Prior to 2012,  
85 local strains contributed to low but increasing sporadic infection rates on the Northeast US coast  
86 (<https://www.cdc.gov/vibrio/surveillance.html>, 2017) (21), with ST631 emerging as the major  
87 lineage that is endemic to near-shore areas of the Atlantic Ocean bordering North America (the  
88 northwest Atlantic Ocean) (31). However, we previously identified a single ST631 isolate

89 lacking hemolysins (21, 27) suggesting this pathogen lineage may have recently evolved through  
90 VP*aI* acquisition.

91         The goal of our study was to understand the genetic events and changing population  
92 context for the evolution of the ST631 pathogenic lineage. We conducted whole and core  
93 genome phylogenetic analysis of three environmental and 39 clinical ST631 isolates along with  
94 isolates from other emergent lineages from the region, which revealed two ST631 clades of  
95 common ancestry, from which human pathogens have evolved in parallel. The single clade I  
96 clinical isolate acquired a *recA* gene insertion previously seen associated with Asian strains, and  
97 had VP*aI* that is typical of isolates harboring *trh* in the absence of *tdh*. In contrast, the clonal  
98 ST631 clade II that dominates Atlantic-derived ST631 infections (31) had a related but distinct  
99 VP*aI*. This VP*aI* contained a *tdh* gene inserted within, not next to, an existing *ure-trh*-T3SS2β  
100 island in close proximity to the *ure* cluster. Nearly all emergent resident lineages and invasive  
101 strains, including all three Pacific lineages (ST36, ST636 and ST43) contained islands that  
102 similarly had a *tdh* gene inserted within the VP*aI* in an identical location adjacent to the *ure*  
103 cluster providing a mechanism for simultaneous acquisition of both hemolysins with T3SS2β.

104

## 105 **RESULTS**

106 **Atlantic endemic ST631 and several invasive strains harboring both the *tdh* and *trh***  
107 **hemolysin genes are clinically prevalent in four reporting Northeast US States.**

108         Ongoing analysis of clinical isolates revealed that even as the Pacific-derived ST36  
109 lineage continued to dominate infections (50%), the endemic (autochthonous) ST631 lineage  
110 accounted for 14% of infections (Table 1). Concurrently, a limited number of other lineages  
111 contributed individually to fewer infections (≤3% each), among which were two lineages that

112 have caused infections in the Pacific Northwest in prior decades: ST43 and ST636 (22, 23).  
113 ST43 and ST636 only recently (2013 and 2011 respectively) (21) have been linked to product  
114 harvested from waters along the Northeast US coast, and also caused infections in subsequent  
115 years. As is common among US clinical isolates, pathogenic strains of all the aforementioned  
116 lineages harbor both the *tdh* and *trh* hemolysin genes (Table 1). Among environmental isolates,  
117 ST34 and ST674 are the most frequently recovered pathogen lineages but these caused  
118 comparatively few infections (Table 1). ST34 was first reported from the environment in 1998,  
119 from both the Gulf of Mexico and near-shore areas of MA, and was also recovered in NH in  
120 2012 (21) suggesting it is an established resident in the region. ST674 which was first reported  
121 from an infection in Virginia in 2007 (32) was first recovered from the local environment in  
122 2012 ([www.pubmlst.org/vparahaemolyticus](http://www.pubmlst.org/vparahaemolyticus)) (21). Notably even though all four ST674  
123 environmental isolates, like ST34, harbored both hemolysin genes, the single ST674 clinical  
124 isolate (MAVP-21) lacked hemolysins (Table 1) (21). The decrease in clinical prevalence of *trh*-  
125 harboring Atlantic endemic ST1127, which caused no infections in the last three years, coincided  
126 with the increase in clinical prevalence of all three Pacific-derived lineages which harbor both  
127 hemolysins. Notably, very few other clinical isolates harbored *trh* in the absence of *tdh* and  
128 clinical isolates containing only *tdh* (i.e. ST1725) were extremely rare (Table 1). Concurrent with  
129 this shift in composition of clinical lineages that includes multiple Pacific-derived strains,  
130 hemolysin producers have increased in relative abundance in nearshore areas of the region,  
131 where historically these represented ~1% of all isolates (27). Since 2012, hemolysin producers  
132 have been recovered more frequently, and in the last two years their proportion has increased by  
133 up to an order of magnitude (comprising as much as 10%) in some regional shellfish associated  
134 populations (data not shown).

135

136 **A single clinical ST631 lineage isolate with an unusual *recA* allele harbors *trh* in the**  
137 **absence of *tdh***

138       Employing ST631-specific marker-based assays (see methods), we identified two  
139 additional 2015 environmental isolates (one from NH and one from MA) and one additional  
140 2011 local-source clinical isolate (MAVP-R) (21) with a hemolysin profile (*trh*<sup>+</sup> without *tdh*)  
141 that is atypical of the ST631 lineage (Table 1). Although analysis of the seven-housekeeping  
142 gene allele combination confirmed the environmental isolates were indeed ST631, MAVP-R was  
143 not ST631 based on only one locus: *recA*. Examination of the *recA* locus of MAVP-R uncovered  
144 a large insertion within the ancestral ST631 *recA* gene (allele *recA*21;  
145 [www.pubmlst.org/vparahemolyticus](http://www.pubmlst.org/vparahemolyticus)) incorporating an intact but different *recA* gene into the  
146 locus [allele *recA*107(33)] and fragmenting the ancestral gene (Fig. 1). The insertion in the  
147 ancestral *recA* gene in MAVP-R is identical to one observed in the *recA* locus of two Hong Kong  
148 isolates (strains S130 and S134) and similar to the one in strain 090-96 (ST189a) isolated in Peru  
149 but believed to have originated in Asia (33).

150

151 **ST631 forms two divergent clades**

152       The existence of three different hemolysin profiles (Table 1) among all available ST631  
153 draft genomes suggested there could be more than one ST631 lineage. Therefore we evaluated  
154 whole genome maximum likelihood (ML) phylogenies of select ST631 isolates and all other  
155 lineages causing two or more infections reported in four Northeast US States to evaluate whether  
156 there was more than one ST631 lineage (Table 1) (Fig. 2). The phylogenetic tree showed that  
157 ST631 isolates, regardless of their hemolysin genotype, clustered together but they formed two

158 distinct clades, indicative of common ancestry (Fig. 2). Clade I harbored either *trh* or no  
159 hemolysins and consisted of all three environmental isolates which were from MA and NH, and  
160 the single clinical isolate MAVP-R, whereas clade II consisted of all other isolates all of which  
161 harbor both hemolysins. The two distinct ST631 clades shared 85% of their DNA in common  
162 and displayed polymorphisms in  $\leq 12\%$  of the shared DNA content. The most closely related  
163 sister lineage to ST631 was formed by *trh*-harboring ST1127 isolates that have been exclusively  
164 reported from clinical sources in the Northeast US (21).

165 We next evaluated the relationships of all available ST631 isolate genomes at NCBI and  
166 sequenced by us (Supplemental Table 1) using a custom core genome multi-locus sequence  
167 typing (cgMLST) method as previously described (31). Minimum spanning trees built from core  
168 genome loci from 42 ST631 isolates indicated that only 390 loci varied between the most closely  
169 related isolate of clade I (MAVP-L) and clade II (G6928) (Fig. 3). The most distantly related  
170 isolates within clade I (G149 and MAVP-R) exhibited 80 core genome loci differences whereas  
171 clade II is clonal with only 51 variant loci between the most divergent isolates: clinical isolate  
172 09-4436 and environmental isolate S487-4, both reported from PEI Canada (Fig. 3) (31).

173

174 **Each ST631 clade independently acquired a distinct pathogenicity island positioned on**  
175 **different chromosomes**

176 Given the variation in ST631, comparisons between these isolates could elucidate the  
177 events that led not only to the evolution of two pathogenic clades but also address unresolved  
178 questions about the unique configurations and contents of pathogenicity islands in western  
179 Atlantic Ocean emergent lineages. The physical proximity of *tdh* with the *ure* cluster and *trh*,  
180 and the co-occurrence of *tdh* with T3SS2 $\beta$  reported in many *tdh*<sup>+</sup>/*trh*<sup>+</sup> clinical isolates suggested



181 *tdh* could be harbored within or next to the same pathogenicity island harboring *trh* in at least  
182 some lineages as was previously suggested (20, 24, 34).

183 To identify the location and determine the architecture of the pathogenicity elements  
184 harboring hemolysin genes, we generated high quality annotated genomes for the clade I ST631  
185 isolate MAVP-R and clade II ST631 isolate MAVP-Q (both reported in 2011 from MA)  
186 employing PacBio sequencing. The pathogenicity island regions in these isolates genomes were  
187 extracted, aligned, and the contents compared with pathogenicity island harboring two *tdh* genes  
188 [previously called Vp-PAI (15), VP $\alpha$ I-7 (4) and *tdh*VP $\alpha$ (17)] from RIMD 2210366 and Vp-  
189 PAI<sub>TH3996</sub> (16) [also called *trh*VPI (17)] harboring *trh* (Supplemental Table 2). This comparison  
190 revealed that MAVP-R harbored a pathogenicity island typical of *trh*-containing strains that  
191 includes a linked *ure* cluster and T3SS2 $\beta$  that is orthologous, with the exception of few unique  
192 regions, with Vp-PAI<sub>TH3996</sub> (16) (Supplemental Table 2 and Fig. 4). Because the lack of  
193 convention in uniformly naming syntenous islands that distinguish them from distinctive and yet  
194 functionally analogous islands can impede communication, we hereafter will consistently  
195 reference the same island by a common descriptive name regardless of strain lineage. Hereafter  
196 we will refer to islands sharing the same general configuration to that in MAVP-R by the name  
197 VP $\alpha$ I $\beta$ , and refer to *tdh*-containing islands similar to that described in strain RIMD 2210366 by  
198 the name VP $\alpha$ I $\alpha$ , regardless of strain background. We adopted this simplified nomenclature in  
199 reference to the version of the key virulence determinant carried in the islands (T3SS2 $\alpha$  and  
200 T3SS2 $\beta$ ) in the two already described island types. This scheme importantly accommodates  
201 naming of additional uniquely-configured islands as they are identified. As noted previously (16,  
202 17, 20), VP $\alpha$ I $\beta$  is dissimilar to VP $\alpha$ I $\alpha$  in most gene content with ~ 78 ORFs unique to VP $\alpha$ I $\beta$   
203 (where the number of identified ORFs used for comparison can differ slightly depending on

204 which annotation program is applied) (Supplemental Table 2, Fig. 4). Even so, VP*α*<sub>β</sub> had many  
205 homologous genes of varying sequence identity (n=~38 ORFs, excluding *tdh* homology with *trh*)  
206 when compared to VP*α*<sub>α</sub> (Supplemental Table 2, Fig. 4)(4, 5, 16). Identification of some  
207 homologs required that we relax matching to 50% such as for the divergent, but homologous  
208 T3SS2*α* and T3SS2*β* genes encoding the apparatus, chaperones, and some shared effectors  
209 (Supplemental Table 2). VP*α*<sub>β</sub> from strain TH3996 and VP*α*<sub>α</sub> from pandemic strain RIMD  
210 2210633 are inserted in an identical location in chromosome II adjacent to an Acyl-CoA  
211 hydrolase-encoding gene. In contrast the VP*α*<sub>β</sub>s in MAVP-R, ST1127 isolate MAVP-25, and  
212 Asia-derived AQ4037 are in chromosome I, in each case in the same insertion location identified  
213 for strain AQ4037 (17).

214       MAVP-Q contained both *tdh* and *trh* within the same contiguous unique VP*α* (hereafter  
215 called VP*α*<sub>γ</sub>) that shared features with both VP*α*<sub>α</sub> and VP*α*<sub>β</sub> (Fig. 4, Supplemental Table 2).  
216 Specifically, VP*α*<sub>γ</sub> had a core that with few exceptions was orthologous in content and  
217 syntenous with VP*α*<sub>β</sub> from MAVP-R (Fig. 4). VP*α*<sub>γ</sub> displays high conservation with VP*α*<sub>α</sub>  
218 near its 3' end, as has been described in other draft *tdh*<sup>+</sup>*trh*<sup>+</sup> harboring genomes (20) as well as in  
219 the VP*α*<sub>β</sub> island of strain TH3996, although the presence of this element may not be typical of  
220 VP*α*<sub>β</sub> (e.g. it is absent in the islands from AQ4037 (17), MAVP-R and MAVP-25). The VP*α*<sub>γ</sub>  
221 also contained a *tdh* gene homologous to *tdh2* (also called *tdhA*) from VP*α*<sub>α</sub> (98.6%) near its 5'  
222 end but not at the 5' terminus of the island (Fig. 4). Rather, the DNA flanking both sides of the  
223 *tdh* gene in VP*α*<sub>γ</sub> was conserved in VP*α*<sub>β</sub> of MAVP-R and absent from VP*α*<sub>α</sub>, (Fig. 4).  
224 Analysis of 300 genomes of *V. parahaemolyticus* (representing a minimum of 28 distinct  
225 sequence types) of sufficient quality for analysis confirmed that the module of four hypothetical  
226 proteins preceding the *tdh2* homolog was present only in *trh*-harboring genomes, but not in

227 genomes harboring *tdh* in the absence of *trh* (i.e. VPα containing genomes), providing  
228 evidence that the *tdh* gene was acquired horizontally by insertion into, not next to, an existing  
229 VPβ, perhaps through activity of the adjacent transposase gene (11) (Supplemental Table 3,  
230 Supplemental fig. 1, and data not shown). Like with VPα from RIMD 2210633, and VPβ of  
231 TH3996, VPγ of clade II ST631 is located in a conserved location of chromosome II, adjacent  
232 to an Acyl-CoA hydrolase-encoding gene.

233         The final environmental ST631 clade I isolate that lacked hemolysins, G149, had no  
234 VPα, β or γ elements in its genome. Close examination of the DNA corresponding to the VPα  
235 insertion sites in either chromosome revealed no remnants of these islands in either chromosomal  
236 location indicating this isolate likely never acquired a pathogenicity island (Supplemental Fig. 2  
237 and data not shown). Because clade I isolate G149 lacked these islands, this could be the  
238 ancestral state of the ST631 lineage (21).

239

240 **Most clinically prevalent strains from the Northeast US harbor similar contiguous**  
241 **pathogenicity islands containing *tdh* inserted in the same location of their VPα**

242         We next asked which other strains likely residing within the mixed population with  
243 ST631 in near-shore areas of the Northeast US harbored islands of similar structure to VPγ that  
244 contain both hemolysin genes. Assembly of short-read sequences into contigs that cover the full  
245 length of VPα which is necessary for comparative analysis of entire island configuration was  
246 impeded by the fact that homologous transposase sequences were repeated multiple times  
247 throughout the island. Therefore, we determine whether other lineages harboring both hemolysin  
248 genes harbor *tdh* in the same island location, between the conserved VPβ/γ module of four  
249 hypothetical proteins (to the left or 5' of *tdh*) and the *ure* cluster (to the right or 3' of *tdh*) (Fig. 4)

250 by combining bioinformatics analysis of sequenced genomes with amplicon assays  
251 (Supplemental Fig. 1). First we analyzed assembled draft genomes for *tdh* co-occurrence and  
252 proximity with the four adjacent hypothetical protein-encoding genes (See Methods). Every  
253 emergent pathogenic lineage (Table 1) harboring both *tdh* and *trh* carried homologous DNA  
254 corresponding to all four hypothetical proteins adjacent to the *tdh* gene in a contiguous segment  
255 (Supplemental Table 3). To determine whether *tdh* was also adjacent to the *ure* cluster in these  
256 same strains we next designed specific flanking primers and amplified the unique juncture  
257 between the *tdh*-containing transposon associated module and the *ure* cluster for all clinical  
258 isolates harboring both *tdh* and *trh* (See Methods) (Supplemental Fig. 1). The results were  
259 congruent with our bioinformatics assessment (Supplemental Table 3), and demonstrated that  
260 isolates from all emergent pathogenic lineages harboring both hemolysins have *tdh* inserted in  
261 close proximity to an *ure* cluster in a configuration similar to VP $\alpha$ I $\gamma$  from MAVP-Q (Fig. 5,  
262 Table 1). This confirmed that these strains harboring both hemolysins harbor *tdh* within, and not  
263 next to, the same VP $\alpha$ I thereby facilitating simultaneous acquisition of both hemolysin genes.

264

## 265 **DISCUSSION**

266 Even preceding the increased illnesses from Pacific-invasive lineages, two different  
267 clades of the predominant endemic Atlantic lineage of pathogenic *V. parahaemolyticus*, ST631  
268 (31) evolved and contributed to a rise in sporadic illnesses in the four reporting Northeast US  
269 States (Table 1, Fig. 2 & 3). Several lines of evidence support the interpretation of parallel  
270 pathogen evolution. The two lineages exhibit differences in both clinical and environmental  
271 prevalence suggesting the pathogenic variants of each clade have not evolved the same degree of  
272 virulence (Table 1). Pathogenic members in each lineage also acquired different pathogenicity

273 islands with different hemolysin gene content (Fig. 2 & 3). Although it was a formal possibility  
274 that ST631 clade II evolved from clade I by independent horizontal acquisition of *tdh* into its  
275 existing VPαIβ, it is notable that other resident and even invasive lineages now in the Atlantic  
276 harbor VPαIγ with *tdh* inserted into the same location, suggesting a common evolutionary origin  
277 of this hybrid type island (Fig. 4 and Supplemental Fig. 1). Finally, each of the two clades harbor  
278 VPαI insertions on different chromosomes: the less clinically prevalent ST631 clade I contains  
279 three isolates that harbor VPαIβ in chromosome I (Fig. 3) and a single environmental isolate  
280 lacking any island (Table 1, supplemental Fig. 2), whereas the clonal ST631 clade II isolates all  
281 harbor VPαIγ on chromosome II.

282         Given that several other resident lineages harbor similar β and γ-type VPαI, pathogens in  
283 each clade could have acquired their islands from the reservoir of strains already circulating in  
284 the Atlantic even before the presume arrival of invasive Pacific lineages. Several well-  
285 documented members of the Gulf of Mexico *V. parahaemolyticus* population (35-37) may also  
286 have expanded their range through movement of ocean currents and could be the source for these  
287 VPαI (Table 1, Fig. 5). But historically, hemolysin producers were extremely rare in near shore  
288 areas of the Atlantic US coast (25) and represented only about ~1% of isolates in an estuary of  
289 NH as of a decade ago (27) limiting the potential for interacting partners or sources for acquired  
290 VPαI. Given this historical context, it is remarkable that two different clades from the same  
291 lineage independently acquired different VPαI-which for clade II ST631 occurred prior to 2007 -  
292 well before the recent shift in abundance of hemolysin producers.

293         The parallel evolution of two different lineages through lateral DNA acquisition alludes  
294 to the possibility that as-yet-undefined attributes may increase the chances of acquisition or  
295 prime some bacterial lineages (such as ST631) to more readily acquire and maintain genetic

296 material or become pathogenic upon island acquisition. Even though the ecological niche in  
297 which horizontal island acquisition took place is unknown, it is conceivable that co-colonization  
298 of hosts or substrates favorable to the growth of ST631 and hemolysin producers may have  
299 facilitated island movement. Certainly, association of bacteria with specific marine substrates  
300 such as chitinous surfaces of plankton that also induce a natural state of competence could  
301 promote lateral transfer through close contact between the progenitors of the pathogenic  
302 subpopulation of each clade and island donors (3, 38, 39). Alternatively, conjugative plasmids or  
303 transducing phage could have been the agents of island delivery. The finding that the only  
304 clinical clade I isolate, MAVP-R, also harbors a second horizontal insertion in its *recA* locus that  
305 matched one previously found in Asia-derived strains (33) indicates it acquired more than one  
306 segment of foreign DNA during its evolution as a pathogen (Fig. 1) further illustrating that  
307 mechanisms that facilitate DNA transfer and acquisition may both have been at play. It also  
308 suggests that horizontal transfer of DNA from introduced strains not yet detected in the Atlantic  
309 could add to the genetic material available for pathogen evolution from Atlantic Ocean  
310 populations. The more detailed molecular epidemiological, comparative genomic, and functional  
311 analyses necessary to assess the impact of introduced pathogens on resident Atlantic lineages are  
312 warranted given this evidence and the documented introduction of multiple Pacific-derived  
313 strains in the region (Table 1).

314         There has been some consideration of the roles of human virulence determinants in  
315 ecological fitness, but the natural context of pathogenic *V. parahaemolyticus* evolution is still  
316 unknown (40-42). Whereas *tdh* and T3SS2 $\alpha$  each may promote growth when bacteria are under  
317 predation, isolates that carry *trh*-containing islands (which likely also have T3SS2 $\beta$ ) do not  
318 derive similar benefits from their islands (43). This is surprising considering the islands encode

319 several homologous effectors (Fig. 4 and Supplemental Table 2) that are not thought to  
320 contribute to human disease but could mediate eukaryotic cell interactions with natural hosts  
321 thereby promoting environmental fitness (13, 14). The general lack of knowledge of unique  
322 T3SS2 $\beta$  effectors and other gene function in these islands (Fig. 4 and Supplemental Table 2)  
323 even with regard to human disease, limits comparative analysis with the well-studied and  
324 functionally defined VP $\alpha$  which could elucidate the bases for pathogen evolution. The higher  
325 clinical prevalence of clade II ST631 than clade I which has also been recovered on more than  
326 one occasion from the environment (Table 1) could indicate that VP $\gamma$  confers greater virulence  
327 potential than VP $\beta$ , perhaps owing to the presence of *tdh*, a known virulence factor (1, 7, 44).  
328 However, the resident community members in both the Pacific and the Atlantic Ocean that  
329 harbor *tdh* and T3SS2 $\alpha$  comparatively rarely cause human infections (21-23). The unique  
330 environmental conditions that underlie pathogen success from northern latitudes that favors  
331 strains with VP $\beta$  and VP $\gamma$  including two different ST631 lineages suggests the shared content  
332 of these islands could confer abilities that are distinct from VP $\alpha$  which could underlie the  
333 repeated acquisition and maintenance of these related islands by so many different lineages now  
334 present in near-shore areas of the Northeast US.

335

## 336 **MATERIALS AND METHODS**

### 337 **Bacteria isolates, media and growth conditions.**

338 *V. parahaemolyticus* clinical isolates for this study were provided by cooperating public  
339 health laboratories in Massachusetts, New Hampshire, Maine, and Connecticut whereas a select  
340 number of environmental isolates were enriched from estuarine substrates as described (21).  
341 Detailed information about these isolates was described previously (31) and listed in

342 Supplemental Table 1. Isolates were routinely cultured in Heart Infusion (HI) media  
343 supplemented with NaCl at 37°C as described (21).

344

345 **Whole genome sequencing, assembly, annotation and sequence type identification.**

346 Genomic DNA was extracted using the Wizard Genomic DNA purification Kit (Promega,  
347 Madison WI USA) or by organic extraction (21). The qualities of all the genomic DNA was  
348 measured by NanoDrop (ThermalFisher, Waltham MA USA). Libraries for DNA sequencing  
349 were prepared using a high-throughput Nextera DNA preparation protocol (45) using an optimal  
350 DNA concentration of 2ng/μl. Genomic DNA was sequenced using an Illumina – HiSeq2500  
351 device at the Hubbard Center for Genome Studies at the University of New Hampshire, using a  
352 150bp paired-end library. *De novo* assembly was performed using the A5 pipeline (46), and the  
353 assemblies annotated with Prokka1.9 using the "genus" option and selecting "*Vibrio*" for the  
354 reference database (47). The sequence types were subsequently determined using the SRST2  
355 pipeline (48). The sequence type of each genome was determined when using *V.*  
356 *parahaemolyticus* as the database (<https://pubmlst.org/vparahaemolyticus/>). For most isolates  
357 where the combination of each allele was not found in the database representing novel sequence  
358 types, the genome was submitted for a new sequence type designation  
359 ([www.pubmlst.org/vparahaemolyticus](http://www.pubmlst.org/vparahaemolyticus)).

360 Isolates MAVP-Q and MAVP-R were sequenced using the Pacific Biosciences RSII  
361 technology. Using between 3.7-5.3 μg DNA, the library preparation and sequencing was  
362 performed according to the manufacturer's instructions (Pacific Biosciences, Menlo Park CA,  
363 USA) and reflects the P5-C3 sequencing enzyme and chemistry for MAVP-Q isolate and the P6-  
364 C4 configuration for MAVP-R. The mass of double-stranded DNA was determined by Qubit



365 (Waltham, MA USA) and the sample diluted to a final concentration of 33  $\mu\text{g} / \mu\text{L}$  in a volume  
366 of 150  $\mu\text{L}$  elution buffer (Qiagen, Germantown MD USA). The DNA was sheared for 60  
367 seconds at 4500 rpm in a G-tube spin column (Covaris, Woburn MA USA) which was  
368 subsequently flipped and re-spun for another 60 seconds at 4500 rpm resulting in a  $\sim 20,000$  bp  
369 DNA verified using a DNA 12000 Bioanalyzer gel chip (Agilent, Santa Clara, CA USA). The  
370 sheared DNA isolate was then re-purified using a 0.45X AMPure XP purification step (Beckman  
371 Coulter, Indianapolis IN USA). The DNA was repaired by incubation in DNA Damage Repair  
372 solution. The library was again purified using 0.45X Ampure XP and SMRTbell adapters ligated  
373 to the ends of the DNA at 25°C overnight. The library was treated with an exonuclease cocktail  
374 (1.81 U/ $\mu\text{L}$  Exo III 18 and 0.18 U/ $\mu\text{L}$  Exo VII) at 37°C for 1 hour to remove un-ligated DNA  
375 fragments. Two additional 0.45X Ampure XP purifications steps were performed to remove  
376  $<2000$  bp molecular weight DNA and organic contaminant.

377           Upon completion of library construction, samples were validated using an Agilent  
378 DNA 12000 gel chip. The isolate library was subjected to additional size selection to the range  
379 of 7,000 bp – 50,000 bp to remove any SMRTbells  $< 5,000$  bp using Sage Science Blue Pippin  
380 0.75% agarose cassettes to maximize the SMRTbell sub-read length for optimal *de*  
381 *novo* assembly. Size-selection was confirmed by Bio-Analysis and the mass was quantified using  
382 the Qubit assay. Primer was then annealed to the library (80°C for 2 minute 30 followed by  
383 decreasing the temperature by 0.1°/s to 25°C). The polymerase-template complex was then  
384 bound to the P5 or P6 enzyme using a ratio of 10:1 polymerase to SMRTbell at 0.5 nM for 4  
385 hours at 30°C and then held at 4°C until ready for magbead loading, prior to sequencing. The  
386 magnetic bead-loading step was conducted at 4°C for 60-minutes per manufacturer's guidelines.  
387 The magbead-loaded, polymerase-bound, SMRTbell libraries were placed onto the RSII machine

388 at a sequencing concentration of 110-150 pM and configured for a 180-minute continuous  
389 sequencing run. Long read assemblies were constructed using HGAP version 2.3.0 for *de novo*  
390 assembly generation. Further, hybrid assemblies were generated and error corrected with  
391 illumina raw reads using Pilon v1.20 (49).

392

### 393 **Lineage-specific marker-based assays**

394 To more rapidly identify ST631 isolates from clinical and environmental collections we  
395 developed PCR-amplicon assays to unique gene content in ST631. Whole genome comparisons  
396 were performed on MAVP-Q (ST631 clinical strain), G149 (ST631 environmental strain),  
397 MAVP-26 (ST36), RIMD2210633 (ST3), and AQ4037 (ST96) (Supplemental Fig. 3). A total of  
398 26 distinct genomic regions, each greater than 1kb in size, were present in MAVP-Q but absent  
399 in other comparator genomes, including environmental ST631 that lacks hemolysins (G149)  
400 (Supplemental Fig. 3). Within a large genomic island ~37.6 Kb in length with an integrase at one  
401 terminus and an overall lower GC content (40.6% compared to 45.8% for the genome) a single  
402 ORF homologous to restriction endonucleases (AB831\_06355) that was restricted to clinical  
403 ST631 isolates in our collection and publicly available draft genomes (n=693)  
404 (<http://www.ncbi.nlm.nih.gov/genome/691>, 2017) was selected as a suitable amplicon target. The  
405 distribution of this locus was further analyzed using the BLAST algorithm by a query against the  
406 nucleotide collection, the non-redundant protein sequences, and against the genus *Vibrio* (taxid:  
407 662), excluding *V. parahaemolyticus* (taxid: 691), using the default settings for BLASTn (50).  
408 Similar approaches were applied to identify ST631 diagnostic loci inclusive of the single  
409 environmental isolate (G149), which identified a hypothetical protein encoding region  
410 (AB831\_06535) (ST631env). Oligonucleotide primers were designed to amplify the diagnostic

411 regions including AB831\_06355 using primers ST631*end* F  
412 (5'AGTTCATCAGGTAGAGAGTTAGAGGA3') and ST631*end*R  
413 (5'TCTTCGTTACCATAGTATGAGCCA3') which produces an amplicon of c.a. 494bp, and  
414 AB831\_06535 using primers ST631*env*F (5'TGGGCGTTAGGCTTTGC3') and ST631-*env*R  
415 (5'GGGCTTCTACGACTTTCTGCT3') producing an amplicon of 497bp.

416       Amplification of diagnostic loci was evaluated in individual assays using genomic DNA  
417 from positive and negative controls: MAVP-Q and G149 (ST631), G4186 (ST34), G3578  
418 (ST674), and MAVP-M (ST1127), MAVP-26 (ST36) and G61 (ST1125). Amplification of  
419 specific sequence types were performed with Accustart enzyme mix on purified DNA. Cycling  
420 was performed with an initial denaturation at 94°C for 3 min., followed by 30 cycles of a  
421 denaturation at 94°C for 1min, annealing at 55°C for 1 min, and amplification at 72°C for 30s  
422 with a final elongation at 72°C for 5 min. The primer pairs only produced amplicons from  
423 template DNA from ST631 and each was the expected size (data not shown, and Supplemental  
424 Fig. 3). Amplicon assays were applied to 208 clinical isolates from the Northeast US States (ME,  
425 NH, MA and CT) and 1140 environmental isolates collected from 2015-2016 from NH and MA.  
426 These assays identified all known ST631 clinical isolates with 100% specificity and also  
427 identified an additional 7 *tdh*<sup>+</sup>*trh*<sup>+</sup> clinical isolates (ST631*end* and ST631*env* positive), and two  
428 environmental (ST631*end* negative and ST631*env* positive) isolates from our archived collection.  
429 Each, with the exception of MAVP-R, was subsequently confirmed to be ST631 by seven-locus  
430 MLST ([www.pubmlst.org](http://www.pubmlst.org)).

431

432 **Examination of *recA* allele and adjacent sequences**

433 The PacBio sequenced genome of MAVP-R, contig 000001 (Accession No.  
434 MPPP00000000) that contained the *recA* gene, was annotated using PROKKA1.9 (47). The  
435 sequences of *recA* and its surrounding DNA was then compared to the contig containing *recA*  
436 region from strain S130 (AWIW01000000), S134 (AWIS01000000), 090-96 (JFFP01000036)  
437 (33) and MAVP-Q (Accession No. MDWT00000000 ). The map of *recA* region of the five  
438 isolates was illustrated using Easyfig (51).

439

#### 440 **Core genome SNP determination and phylogenetic analysis**

441 Whole genome phylogenies were constructed with single nucleotide polymorphisms  
442 (SNPs) identified from draft genomes using kSNP3 to produce aligned SNPs in FASTA format  
443 (52). A maximum likelihood (ML) tree was then built from the FASTA file using raxMLHPC  
444 with model GTRGAMMA, -f a and 100 bootstraps (53).

445 Minimum spanning tree (MST) analysis was built based on core gene SNPs produced  
446 from a cluster analysis. The cluster analysis of ST631 was performed using a custom core  
447 genome multi-locus sequence type (cgMLST) analysis using RidomSeqSphere+software v3.2.1  
448 (<http://www.ridom.de.seqsphere>, Ridom GmbH, Münster, Germany) as previously described  
449 (31). Briefly, the software first defines a cgMLST scheme using the target definer tool with  
450 default settings using the PacBio generated MAVP-Q genome as the reference. Then, five other  
451 *V. parahaemolyticus* genomes (BB22OP, CDC\_K4557, FDA\_R31, RIMD2210633, and UCM-  
452 V493) were used for comparison with the reference genome to establish the core and accessory  
453 genome genes. Genes that are repeated in more than one copy in any of the six genomes were  
454 removed from the analysis. Subsequently, a task template was created that contains both core and  
455 accessory genes. Each individual gene locus from MAVP-Q was assigned allele number 1. Then

456 each ST631 isolate genome assembly was queried against the task template, where any locus that  
457 differed from the reference genome or any other queried genome was assigned a new allele  
458 number. The cgMLST performed a gene-by-gene analysis of all core genes (excluding accessory  
459 genes) and identified SNPs within different alleles to establish genetic distance calculations.

460

## 461 **Configuration and distribution of VPAs**

462 The VPAl sequence from the PacBio sequenced genomes of MAVP-Q and MAVP-R  
463 were identified by comparison with the published RIMD2210633 VPAl-7 (NC\_004605 region  
464 between VPA1312 – VPA1395) and VPAl<sub>TH3996</sub> (AB455531) (16). Identification of the complete  
465 MAVP-Q VPAl $\gamma$  and genomic junctures in chromosome II was done by comparison with the  
466 same region of chromosome II in MAVP-R and G149 (which lack an island in this location)  
467 using Mauve (54). In a reciprocal manner, the absence of an island in chromosome I in MAVP-Q  
468 and G149 was assessed by comparison with chromosome I of MAVP-R. MAVP-Q VPAl $\gamma$   
469 (MF066646) and MAVP-R VPAl $\beta$  (MF066647) were then extracted as a single contiguous  
470 sequence and annotated using Prokka 1.9. Gene content and order of the VPAl elements in  
471 MAVP-Q, MAVP-R and RIMD2210633 were then illustrated by Easyfig (51). Roary (55) was  
472 then employed to determine homologs among VPAs based on each island's annotated sequences  
473 with identity set at 50%. Identification of the genome locations of VPAl $\beta$  in ST1127 isolate  
474 MAVP-M (accession number GCA\_001023155) and for VPAl $\gamma$  in AQ4037 (accession number  
475 GCA\_000182365) (17) was also done using Mauve (54).

476 To examine the distribution of the VPAl $\gamma$  in all publicly available draft genomes  
477 (<https://www.ncbi.nlm.nih.gov/genome/genomes/691>, 2016) and genomes from archived  
478 regional isolates, whole draft genome sequences were aligned to a 6,118 bp subsequence of the

479 MAVP-Q VP*α*I with NASP version 1.0.2 (56) (<https://pypi.python.org/pypi/nasp/1.0.2>, 2017).  
480 This subsequence spanned the unique juncture of the four conserved hypothetical proteins  
481 (AB831\_22090, AB831\_22095, AB831\_22100, AB831\_22105) with the adjacent inserted *tdh*  
482 (AB831\_22110, c.a. 2549 bp upstream of *ure* cluster)(Supplemental Fig. 1). Percent coverage of  
483 the reference sequence was used to determine whether each genome harbored only the four  
484 hypothetical proteins, only a *tdh* gene, or the entire module including the fusion of the four genes  
485 with *tdh* (Supplemental Fig. 1 and Supplemental Table 3). The sequence type of each genome  
486 harboring the fused element characteristic of VP*α*I $\gamma$  was then determined using the SRST2  
487 pipeline (48). Where sequencing reads were not available as the input for SRST2, they were  
488 simulated from assemblies using an in-house Python script  
489 (<https://github.com/kpdrees/fasta2reads>).

490 A PCR amplification approach was developed and applied to survey the presence of *tdh*  
491 adjacent to the *ure* gene cluster. Primers were designed to conserved sequences of the 3' end of  
492 *tdh* (PIHybF8: 5'GCCAACATGGATATAAATAAAAATGA3') and the 5' end of *ureG*  
493 (*tdhUreGrev*5: 5'GACAAAGGTATGCTGCCAAAAGTG3') as determined by gene alignments,  
494 which when used together produced a 2631 bp amplicon of the insertion juncture when used with  
495 MAVP-Q as a template (Supplemental Fig. 4). Amplification was performed on purified DNA  
496 with Accustart enzyme mix, with an initial denaturation at 94°C for 3 min., followed by 30  
497 cycles of a denaturation at 94°C for 1 min, annealing at 61°C for 1min, and amplification at 72°C  
498 for 2.5 min, with a final elongation at 72°C for 5 min. This amplification was performed in  
499 parallel with a diagnostic multiplex PCR amplification of *tdh*, *trh* and *tlh* using published  
500 methods (10, 57) to investigate the co-occurrence of VP*α*I $\gamma$  with both hemolysin encoding genes

501 in representative isolates of various clinically prevalent sequence types. Amplicons were  
502 visualized using a 1.2% agarose gel in TAE buffer (Supplemental Fig. 4).

503

#### 504 **Nucleotide sequence accession numbers.**

505 The accession number of Pacific Biosciences sequenced genome for MAVP-Q is  
506 MDWT000000000, and for MAVP-R is MPPP000000000. The accession number of Illumina  
507 sequenced draft genome for G6928 is MPPN000000000, for MA561 is MPPM000000000 and for  
508 G149 is MPPO000000000. Detailed information about all other ST631 isolate draft genomes were  
509 described previously (31) and are listed in Supplemental Table 1. The accessions for the short  
510 reads for the remaining sequenced genomes are listed in Supplemental Table 4. The accession  
511 number of VP $\alpha$ I $\beta$  from MAVP-R is MF066647 and the accession number of VP $\alpha$ I $\gamma$  from MAVP-  
512 Q is MF066646.

513

#### 514 **ACKNOWLEDGEMENTS**

515

516 We are grateful for clinical strains and wish to thank specifically: Jana Ferguson and Tracy Stiles  
517 of the Massachusetts Department of Public Health, and M. Hickey and C. Schillaci from the  
518 Massachusetts Department of Marine Fisheries; J.K. Kanwit of the Maine Department of Marine  
519 Resources and A. Robbins from the Maine Department of Health and Human Services; and  
520 Lurn Mank from the Connecticut Department of Public Health Laboratory, and K. DeRosia-  
521 Banick, Connecticut Department of Agriculture, Bureau of Aquaculture. Assistance with genome  
522 sequencing was provided by W. K. Thomas, and technical assistance provided by Jacqueline  
523 Lemaire, Kari Hartman, Christopher Hallee, Michael Malanga, Saba Ilyas, Jeffrey Hall, Joseph

524 Sevigny, Marcus Dillon, Kenneth Flynn, Alison Goupil, Sarah Eggert, Jillian Means, Randi  
525 Foxall, and M. Sabrina Pankey. Partial funding for this work was provided by the USDA  
526 National Institute of Food and Agriculture (Hatch projects NH00574, NH00609 [accession  
527 number 233555], and NH00625 [accession number 1004199]). Additional funding was provided  
528 by the National Oceanic and Atmospheric Administration College Sea Grant program and grants  
529 R/CE-137, R/SSS-2, and R/HCE-3. Support was also provided through the National Institutes of  
530 Health (1R03AI081102-01), the National Science Foundation (EPSCoR IIA-1330641), and the  
531 National Science Foundation (DBI 1229361 NSF MRI). N.G.-E. was funded through the FDA  
532 Foods Science and Research Intramural Program. Feng Xu and Cheryl A. Whistler declare a  
533 potential conflict of interest in the form of a pending patent application (U.S. patent application  
534 62/128,764). This is Scientific Contribution Number 2722 for the New Hampshire Agricultural  
535 Experiment Station.

536

## 537 REFERENCES

- 538 1. **Hiyoshi H, Kodama T, Iida T, Honda T.** 2010. Contribution of *Vibrio*  
539 *parahaemolyticus* virulence factors to cytotoxicity, enterotoxicity, and lethality in mice.  
540 *Infect Immun* **78**:1772-1780.
- 541 2. **Scallan E, Hoekstra RM, Angulo FJ, Tauxe RV, Widdowson M-A, Roy SL, Jones**  
542 **JL, Griffin PM.** 2011. Foodborne illness acquired in the United States—major  
543 pathogens. *Emerg Infect Dis* **17(1)**:7-15.
- 544 3. **Hazen TH, Pan L, Gu J-D, Sobecky PA.** 2010. The contribution of mobile genetic  
545 elements to the evolution and ecology of *Vibrios*. *FEMS Microbiol Ecol* **74**:485-499.



- 546 4. **Hurley CC, Quirke A, Reen FJ, Boyd EF.** 2006. Four genomic islands that mark post-  
547 1995 pandemic *Vibrio parahaemolyticus* isolates. *BMC Genomics* **7**:104  
548 DOI:110.1186/1471-2164-1187-1104.
- 549 5. **Boyd EF, Cohen AL, Naughton LM, Ussery DW, Binnewies TT, Stine OC, Parent**  
550 **MA.** 2008. Molecular analysis of the emergence of pandemic *Vibrio parahaemolyticus*.  
551 *BMC Microbiol* **8**:110.
- 552 6. **Kishishita M, Matsuoka N, Kumagai K, Yamasaki S, Takeda Y, Nishibuchi M.** 1992.  
553 Sequence variation in the thermostable direct hemolysin-related hemolysin (*trh*) gene of  
554 *Vibrio parahaemolyticus*. *Appl Environ Microbiol* **58**:2449-2457.
- 555 7. **Honda T, Ni Y, Miwatani T, Adachi T, Kim J.** 1992. The thermostable direct  
556 hemolysin of *Vibrio parahaemolyticus* is a pore-forming toxin. *Can J Microbiol* **38**:1175-  
557 1180.
- 558 8. **Park K-S, Ono T, Rokuda M, Jang M-H, Iida T, Honda T.** 2004. Cytotoxicity and  
559 enterotoxicity of the thermostable direct hemolysin-deletion mutants of *Vibrio*  
560 *parahaemolyticus*. *Microbiol Immunol* **48**:313-318.
- 561 9. **Shirai H, Ito H, Hirayama T, Nakamoto Y, Nakabayashi N, Kumagai K, Takeda Y,**  
562 **Nishibuchi M.** 1990. Molecular epidemiologic evidence for association of thermostable  
563 direct hemolysin (TDH) and TDH-related hemolysin of *Vibrio parahaemolyticus* with  
564 gastroenteritis. *Infect Immun* **58**:3568-3573.
- 565 10. **Panicker G, Call DR, Krug MJ, Bej AK.** 2004. Detection of pathogenic *Vibrio* spp. in  
566 shellfish by using multiplex PCR and DNA microarrays. *Appl Environ Microbiol*  
567 **70**:7436-7444.

- 568 11. **Nishibuchi M, Kaper JB.** 1995. Thermostable direct hemolysin gene of *Vibrio*  
569 *parahaemolyticus*: a virulence gene acquired by a marine bacterium. Infect Immun  
570 **63**:2093.
- 571 12. **Park K-S, Ono T, Rokuda M, Jang M-H, Okada K, Iida T, Honda T.** 2004.  
572 Functional characterization of two type III secretion systems of *Vibrio parahaemolyticus*.  
573 Infect Immun **72**:6659-6665.
- 574 13. **Broberg CA, Calder TJ, Orth K.** 2011. *Vibrio parahaemolyticus* cell biology and  
575 pathogenicity determinants. Microb Infect **13**:992-1001.
- 576 14. **Zhang L, Orth K.** 2013. Virulence determinants for *Vibrio parahaemolyticus* infection.  
577 Curr Opin Microbiol **16**:70-77.
- 578 15. **Makino K, Oshima K, Kurokawa K, Yokoyama K, Uda T, Tagomori K, Iijima Y,**  
579 **Najima M, Nakano M, Yamashita A.** 2003. Genome sequence of *Vibrio*  
580 *parahaemolyticus*: a pathogenic mechanism distinct from that of *V. cholerae*. The Lancet  
581 **361**:743-749.
- 582 16. **Okada N, Iida T, Park K-S, Goto N, Yasunaga T, Hiyoshi H, Matsuda S, Kodama T,**  
583 **Honda T.** 2009. Identification and characterization of a novel type III secretion system in  
584 trh-positive *Vibrio parahaemolyticus* strain TH3996 reveal genetic lineage and diversity  
585 of pathogenic machinery beyond the species level. Infect Immun **77**:904-913.
- 586 17. **Chen Y, Stine OC, Badger JH, Gil AI, Nair GB, Nishibuchi M, Fouts DE.** 2011.  
587 Comparative genomic analysis of *Vibrio parahaemolyticus*: serotype conversion and  
588 virulence. BMC Genomics **12**:1.
- 589 18. **Zhou X, Gewurz BE, Ritchie JM, Takasaki K, Greenfeld H, Kieff E, Davis BM,**  
590 **Waldor MK.** 2013. *vopZ* A *Vibrio parahaemolyticus* T3SS effector mediates

- 591 pathogenesis by independently enabling intestinal colonization and inhibiting TAK1  
592 activation. Cell Reports **3**:1690-1702.
- 593 19. **Hubbard TP, Chao MC, Abel S, Blondel CJ, zur Wiesch PA, Zhou X, Davis BM,**  
594 **Waldor MK.** 2016. Genetic analysis of *Vibrio parahaemolyticus* intestinal colonization.  
595 Proc Nat Acad Sci USA **113**:6283-6288.
- 596 20. **Ronholm J, Petronella N, Leung CC, Pightling A, Banerjee S.** 2016. Genomic  
597 Features of Environmental and Clinical *Vibrio parahaemolyticus* Isolates Lacking  
598 Recognized Virulence Factors Are Dissimilar. Appl Environ Microbiol **82**:1102-1113.
- 599 21. **Xu F, Ilyas S, Hall JA, Jones SH, Cooper VS, Whistler CA.** 2015. Genetic  
600 characterization of clinical and environmental *Vibrio parahaemolyticus* from the  
601 Northeast USA reveals emerging resident and non-indigenous pathogen lineages. Name:  
602 Front Microbiol **6**:272.
- 603 22. **Banerjee SK, Kearney AK, Nadon CA, Peterson C-L, Tyler K, Bakouche L, Clark**  
604 **CG, Hoang L, Gilmour MW, Farber JM.** 2014. Phenotypic and genotypic  
605 characterization of Canadian clinical isolates of *Vibrio parahaemolyticus* collected from  
606 2000 to 2009. J Clin Microbiol **52**:1081-1088.
- 607 23. **Turner JW, Paranjpye RN, Landis ED, Biryukov SV, González-Escalona N, Nilsson**  
608 **WB, Strom MS.** 2013. Population structure of clinical and environmental *Vibrio*  
609 *parahaemolyticus* from the Pacific Northwest coast of the United States. PLoS ONE  
610 **8(2)**:e55726
- 611 24. **Jones JL, Lüdeke CH, Bowers JC, Garrett N, Fischer M, Parsons MB, Bopp CA,**  
612 **DePaola A.** 2012. Biochemical, serological, and virulence characterization of clinical and  
613 oyster *Vibrio parahaemolyticus* isolates. J Clin Microbiol **50(7)**:2343-2352.

- 614 25. **DePaola A, Ulaszek J, Kaysner CA, Tenge BJ, Nordstrom JL, Wells J, Puhr N,**  
615 **Gendel SM.** 2003. Molecular, serological, and virulence characteristics of *Vibrio*  
616 *parahaemolyticus* isolated from environmental, food, and clinical sources in North  
617 America and Asia. *Appl Environ Microbiol* **69**:3999-4005.
- 618 26. **Haendiges J, Timme R, Allard MW, Myers RA, Brown EW, Gonzalez-Escalona N.**  
619 2015. Characterization of *Vibrio parahaemolyticus* clinical strains from Maryland (2012–  
620 2013) and comparisons to a locally and globally diverse *V. parahaemolyticus* strains by  
621 whole-genome sequence analysis. *Front Microbiol* **6**:125
- 622 27. **Ellis CN, Schuster BM, Striplin MJ, Jones SH, Whistler CA, Cooper VS.** 2012.  
623 Influence of seasonality on the genetic diversity of *Vibrio parahaemolyticus* in New  
624 Hampshire shellfish waters as determined by multilocus sequence analysis. *Appl Environ*  
625 *Microbiol* **78**:3778-3782.
- 626 28. **Nair GB, Ramamurthy T, Bhattacharya SK, Dutta B, Takeda Y, Sack DA.** 2007.  
627 Global dissemination of *Vibrio parahaemolyticus* serotype O3: K6 and its serovariants.  
628 *Clin Microbiol Rev* **20**:39-48.
- 629 29. **Martinez-Urtaza J, Baker-Austin C, Jones JL, Newton AE, Gonzalez-Aviles GD,**  
630 **DePaola A.** 2013. Spread of Pacific Northwest *Vibrio parahaemolyticus* strain. *N Engl J*  
631 *Med* **369**:1573-1574.
- 632 30. **Newton AE, Garrett N, Stroika SG, Halpin JL, Turnsek M, Mody RK, Division of**  
633 **Foodborne W, Environmental D.** 2014. Notes from the field: Increase in *Vibrio*  
634 *parahaemolyticus* infections associated with consumption of Atlantic coast shellfish—  
635 2013. *MMWR Morb Mortal Wkly Rep* **63**:335-336.

- 636 31. **Xu F, Gonzalez-Escalona N, Haendiges J, Myers RA, Ferguson J, Stiles T, Hickey E,**  
637 **Moore M, Hickey JM, Schillaci C.** 2017. Sequence type 631 *Vibrio parahaemolyticus*,  
638 an emerging foodborne pathogen in North America. *J Clin Microbiol* **55**:645-648.
- 639 32. **Lüdeke CH, Gonzalez-Escalona N, Fischer M, Jones JL.** 2015. Examination of  
640 clinical and environmental *Vibrio parahaemolyticus* isolates by multi-locus sequence  
641 typing (MLST) and multiple-locus variable-number tandem-repeat analysis (MLVA).  
642 *Frontiers in microbiology* **6**:564
- 643 33. **González-Escalona N, Gavilan RG, Brown EW, Martinez-Urtaza J.** 2015.  
644 Transoceanic spreading of pathogenic strains of *Vibrio parahaemolyticus* with distinctive  
645 genetic signatures in the *recA* gene. *PloS one* **10**:e0117485.
- 646 34. **Park K-S, Suthienkul O, Kozawa J, Yamaichi Y, Yamamoto K, Honda T.** 1998.  
647 Close proximity of the *tdh*, *trh* and *ure* genes on the chromosome of *Vibrio*  
648 *parahaemolyticus*. *Microbiology* **144**:2517-2523.
- 649 35. **Johnson C, Flowers A, Young V, Gonzalez-Escalona N, DePaola A, Noriega III N,**  
650 **Grimes D.** 2009. Genetic relatedness among *tdh*<sup>+</sup> and *trh*<sup>+</sup> *Vibrio parahaemolyticus*  
651 cultured from Gulf of Mexico oysters (*Crassostrea virginica*) and surrounding water and  
652 sediment. *Microb Ecol* **57**:437-443.
- 653 36. **González-Escalona N, Martinez-Urtaza J, Romero J, Espejo RT, Jaykus L-A,**  
654 **DePaola A.** 2008. Determination of molecular phylogenetics of *Vibrio parahaemolyticus*  
655 strains by multilocus sequence typing. *J Bacteriol* **190**:2831-2840.
- 656 37. **Ellingsen BA, Olsen JS, Granum PE, Rorvik LM, González-Escalona N.** 2013.  
657 Genetic characterization of *trh* positive *Vibrio* spp. isolated from Norway. *Front Cell*  
658 *Infect Microbiol* **3**:107.

- 659 38. **Chen Y, Dai J, Morris JG, Johnson JA.** 2010. Genetic analysis of the capsule  
660 polysaccharide (K antigen) and exopolysaccharide genes in pandemic *Vibrio*  
661 *parahaemolyticus* O3: K6. BMC Microbiol **10**:1.
- 662 39. **Meibom KL, Blokesch M, Dolganov NA, Wu C-Y, Schoolnik GK.** 2005. Chitin  
663 induces natural competence in *Vibrio cholerae*. Science **310**:1824-1827.
- 664 40. **Takemura AF, Chien DM, Polz MF.** 2014. Associations and dynamics of *Vibrionaceae*  
665 in the environment, from the genus to the population level. Front Microbiol **5**:38.
- 666 41. **Lovell CR.** 2017. Ecological fitness and virulence features of *Vibrio parahaemolyticus* in  
667 estuarine environments. Appl Microbiol Biotechnol **101**:1781-1794.
- 668 42. **Johnson CN.** 2013. Fitness factors in vibrios: a mini-review. Microb Ecol **65**:826-851.
- 669 43. **Matz C, Nouri B, McCarter L, Martinez-Urtaza J.** 2011. Acquired type III secretion  
670 system determines environmental fitness of epidemic *Vibrio parahaemolyticus* in the  
671 interaction with bacterivorous protists. PloS one **6**:e20275.
- 672 44. **Nishibuchi M, Kaper JB.** 1985. Nucleotide sequence of the thermostable direct  
673 hemolysin gene of *Vibrio parahaemolyticus*. J Bacteriol **162**:558-564.
- 674 45. **Baym M, Kryazhimskiy S, Lieberman TD, Chung H, Desai MM, Kishony R.** 2015.  
675 Inexpensive multiplexed library preparation for megabase-sized genomes. PloS one  
676 **10**:e0128036.
- 677 46. **Tritt A, Eisen JA, Facciotti MT, Darling AE.** 2012. A5. An integrated pipeline for *de*  
678 *novo* assembly of microbial genomes. PLoS ONE **7**:e42304.
- 679 47. **Seemann T.** 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics.  
680 **30**:2068-9

- 681 48. **Inouye M, Conway TC, Zobel J, Holt KE.** 2012. Short read sequence typing (SRST):  
682 multi-locus sequence types from short reads. *BMC Genomics* **13**:338.
- 683 49. **Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA,**  
684 **Zeng Q, Wortman J, Young SK.** 2014. Pilon: an integrated tool for comprehensive  
685 microbial variant detection and genome assembly improvement. *PloS one* **9**:e112963.
- 686 50. **Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden**  
687 **TL.** 2009. BLAST+: architecture and applications. *BMC Bioinformatics* **10**:421.
- 688 51. **Sullivan MJ, Petty NK, Beatson SA.** 2011. Easyfig: a genome comparison visualizer.  
689 *Bioinformatics* **27**:1009-1010.
- 690 52. **Gardner SN, Slezak T, Hall BG.** 2015. kSNP3. 0: SNP detection and phylogenetic  
691 analysis of genomes without genome alignment or reference genome. *Bioinformatics*  
692 **31**:2877-8.
- 693 53. **Stamatakis A.** 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic  
694 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**:2688-2690.
- 695 54. **Darling AC, Mau B, Blattner FR, Perna NT.** 2004. Mauve: multiple alignment of  
696 conserved genomic sequence with rearrangements. *Genome Res* **14**:1394-1403.
- 697 55. **Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M,**  
698 **Falush D, Keane JA, Parkhill J.** 2015. Roary: rapid large-scale prokaryote pan genome  
699 analysis. *Bioinformatics* **31**:3691-3693.
- 700 56. **Sahl JW, Lemmer D, Travis J, Schupp J, Gillece J, Aziz M, Driebe E, Drees K,**  
701 **Hicks N, Williamson C.** 2016. The Northern Arizona SNP Pipeline (NASP): accurate,  
702 flexible, and rapid identification of SNPs in WGS datasets. *Microb Genom.* **2**:e000074

- 703 57. **Whistler CA, Hall JA, Xu F, Ilyas S, Siwakoti P, Cooper VS, Jones SH.** 2015. Use of  
704 Whole-Genome Phylogeny and Comparisons for Development of a Multiplex PCR Assay  
705 To Identify Sequence Type 36 *Vibrio parahaemolyticus*. J Clin Microbiol **53**:1864-1872.
- 706 58. **Jolley KA, Chan M-S, Maiden MC.** 2004. mlstdbNet—distributed multi-locus sequence  
707 typing (MLST) databases. BMC Bioinformatics **5**:86.
- 708 59. **Alikhan N-F, Petty NK, Zakour NLB, Beatson SA.** 2011. BLAST Ring Image  
709 Generator (BRIG): simple prokaryote genome comparisons. BMC Genomics **12**:402  
710  
711



712 Table 1: Clinical and environmental prevalence of emergent Northeast US *V. parahaemolyticus*  
 713 lineages with associated virulence features.

Sequence type <sup>a</sup>	Northeast US States <sup>b</sup>		MLST Database <sup>c</sup>		Hemolysin genotype	VPaI type <sup>d</sup>
	Clinical	Environmental	Clinical	Environmental		
3	2	0	217	33	<i>tdh</i> <sup>+</sup>	α
36	91	1	58	5	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
631	24	0	12	0	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
	1 <sup>c</sup>	2	0	0	<i>trh</i> <sup>+</sup>	β
	0	1	0	0	neither	absent
43	5	0	17	4	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
636	4	0	2	0	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
1127	4	0	0	0	<i>trh</i> <sup>+</sup>	β
110	3	0	0	1	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
34/324	2	2	4	19	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
674	0	4	1	20	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
	1	0	0	0	neither	absent
308	2	0	0	2	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
12	2	0	0	4	<i>trh</i> <sup>+</sup>	β
162	2	0	1	1	neither	absent
194	2	0	1	0	neither	absent
809	2	0	0	1	<i>trh</i> <sup>+</sup>	β
1716	2	0	0	0	<i>trh</i> <sup>+</sup>	β
1123	1	1	0	0	<i>trh</i> <sup>+</sup>	β
8	1	0	13	5	<i>trh</i> <sup>+</sup>	β
23	1	0	0	3	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
749	1	0	1	0	<i>tdh</i> <sup>+</sup> <i>trh</i> <sup>+</sup>	γ
1295	1	0	0	1	neither	absent
134	1	0	1	0	neither	absent
741	1	0	0	1	neither	absent
98	1	0	0	1	<i>trh</i> <sup>+</sup>	β
1205	1	0	0	1	neither	absent
1561	1	0	0	0	neither	absent
1717	1	0	0	0	neither	absent
1725	1	0	0	0	<i>tdh</i> <sup>+</sup>	α

714 <sup>a</sup> Some clinical isolates had insufficient sequencing coverage to determine sequence type and included  
 715 eight *tdh*<sup>+</sup>*trh*<sup>+</sup> isolates, one *tdh*<sup>+</sup> isolate, four *trh*<sup>+</sup> isolates, and 11 isolates without hemolysins, some of  
 716 which were from wound infections. Two wound infection isolates lacking hemolysins were of known  
 717 sequence types and are not listed above.

718 <sup>b</sup>Data generated from all available gastric infection clinical and environmental isolates four reporting  
 719 Northeast US States including ME, NH, MA, and CT between 2010 and 2016.

720 <sup>c</sup><http://pubmlst.org/vparahaemolyticus>, 2017 (36, 58)

721 <sup>d</sup>Presence of the VPαI architecture was determined by PacBio genome sequencing of strain MAVP-Q  
 722 and MAVP-26, whereas for other strains, identification of VPαI type was determined through illumina  
 723 genome sequencing, PCR amplification and Sanger sequencing.

724 <sup>e</sup>This single isolate harbors a *recA* allele (allele 21) typical of ST631 fused to allele 107 through an  
 725 insertion event, generating a hybrid allele previously described (33).

726

727

728 Figure 1. Schematic of a horizontally acquired insertion in the *recA*-encoding region of MAVP-R.  
729 Sequences of the *recA* gene and flanking region from MAVP-Q (reference ST631 genome),  
730 MAVP-R, Asia-derived isolates S130/S134 and Peru-derived isolate 090-96 were extracted and  
731 aligned. Open reading frames designated with arrows and illustrated by representative colors to  
732 highlight homologous and unique genes. The % similarity between homologs is illustrated by  
733 grey bars.

734

735 Figure 2. Phylogenetic relationships of *V. parahaemolyticus* lineages and identification of  
736 distinct ST631 clades. An ML phylogeny of representative *V. parahaemolyticus* genomes of  
737 clinical strains causing two or more infections was built on whole genome SNPs identified by  
738 reference-free comparisons as described in the methods. The branch length represents the  
739 number of nucleotide substitutions per site. Numbers at nodes represent percent bootstrap  
740 support where unlabeled nodes had bootstraps of less than 70.

741

742 Figure 3. Minimum spanning tree relationships among clade I and clade II ST631. A cgMLST  
743 core gene-by-gene analysis (excluding accessory genes) was performed and SNPs were  
744 identified within different alleles. The numbers above the connected lines (not to scale) represent  
745 SNP differences. The isolates are colored based on different hemolysin genotypes as labeled.

746

747 Figure 4. Comparisons of the pathogenicity islands containing hemolysins and T3SS2.  
748 Sequences of VP*a*I were extracted from select genomes and aligned. VP*a*I $\alpha$  was derived from  
749 ST3 strain RIMD2210633, VP*a*I $\gamma$  was derived from ST631 clade II isolate MAVP-Q, and VP*a*I $\beta$   
750 was derived from ST631 clade I isolate MAVP-R. ORFs are depicted in defined colors and

751 similarities ( $\geq 75\%$ ) among ORFs are illustrated in grey blocks. Homologs between VP $\alpha$  and  
752 VP $\beta/\gamma$  ( $50 > 75\%$  identity) are named and listed in Supplemental Table 2.

753

754 Figure 5. Distribution of VP $\gamma$  in emergent pathogen lineages. The presence of *tdh*, *trh* and  
755 VP $\gamma$  along with positive control *tlh* was determined by PCR amplification using gene-specific  
756 primers and visualized on a 1.2% agarose gel. The order from left to right is 1kb+ ladder, ST3  
757 (MAVP-C), ST36 (MAVP-26), ST631 CII (clade II isolate MAVP-Q), ST631 CI (clade I  
758 isolates MAVP-R and G149), ST43 (MAVP-71), ST636 (MAVP-50), ST1127 (MAVP-M),  
759 ST110 (MAVP-46), ST34 (CTVP19C), ST324 (MAVP-14), and ST674 (CT4291, MAVP21).  
760 The corresponding sizes of the ladder fragments are as labeled to the left and the identity of the  
761 amplicons listed to the right of the gel image.

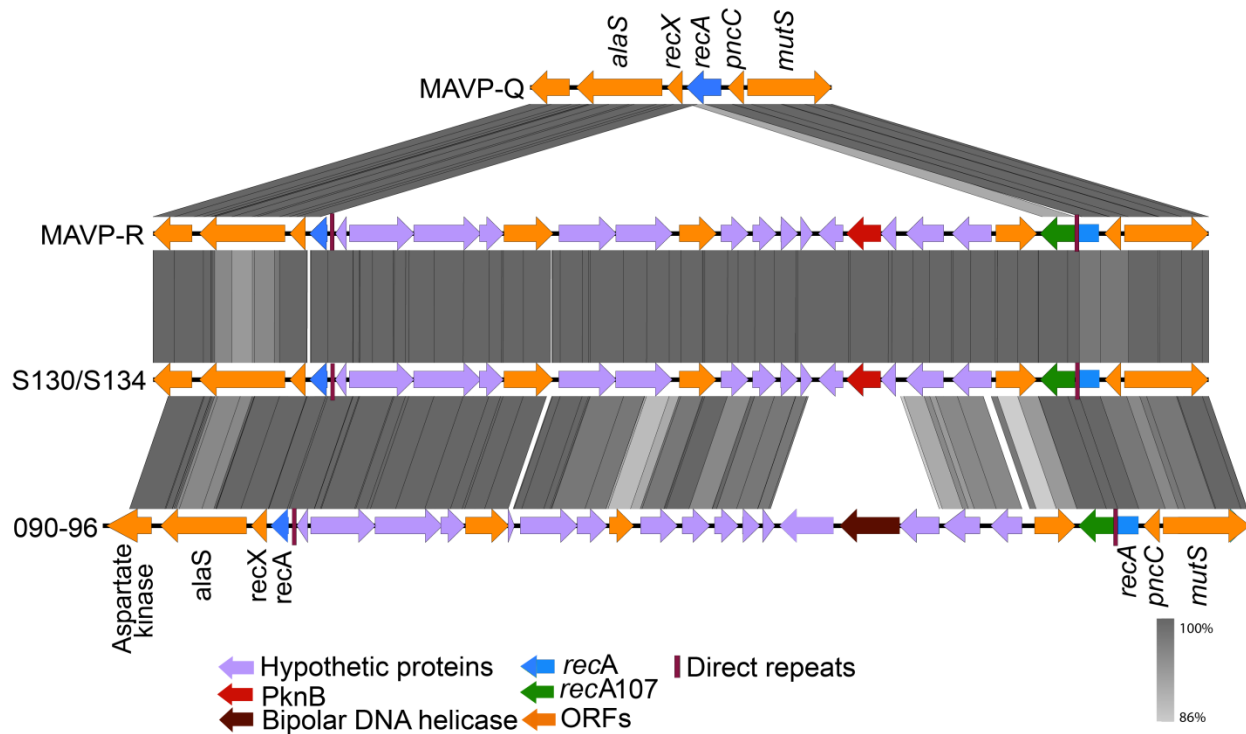


Figure 1. Schematic of a horizontally acquired insertion in the *recA*-encoding region of MAVP-R. Sequences of the *recA* gene and flanking region from MAVP-Q (reference ST631 genome), MAVP-R, Asia-derived isolates S130/S134 and Peru-derived isolate 090-96 were extracted and aligned. Open reading frames designated with arrows and illustrated by representative colors to highlight homologous and unique genes. The % similarity between homologs is illustrated by grey bars.



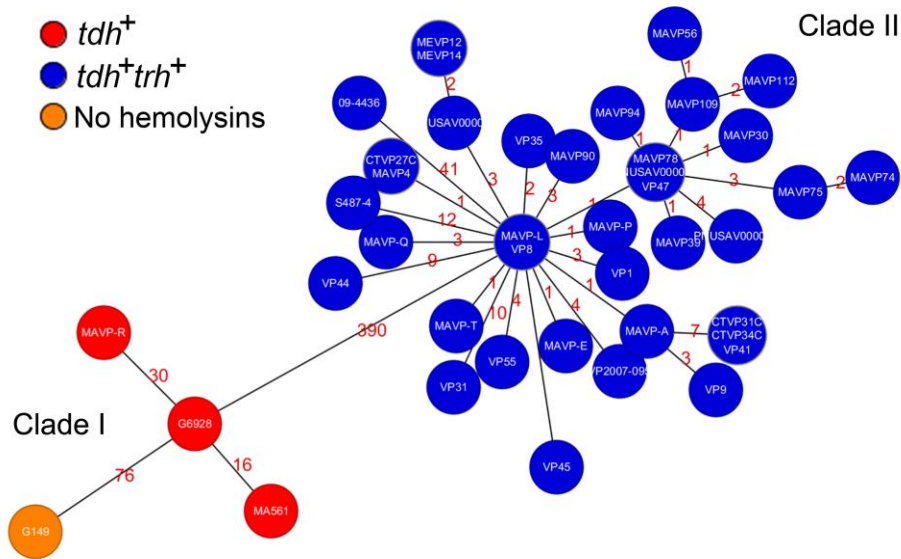


Figure 3. Minimum spanning tree relationships among clade I and clade II ST631. A core gene-by-gene analysis (excluding accessory genes) was performed and SNPs were identified within different alleles. The numbers above the connected lines (not to scale) represent SNP differences. The isolates are colored based on different hemolysin genotypes as labeled.

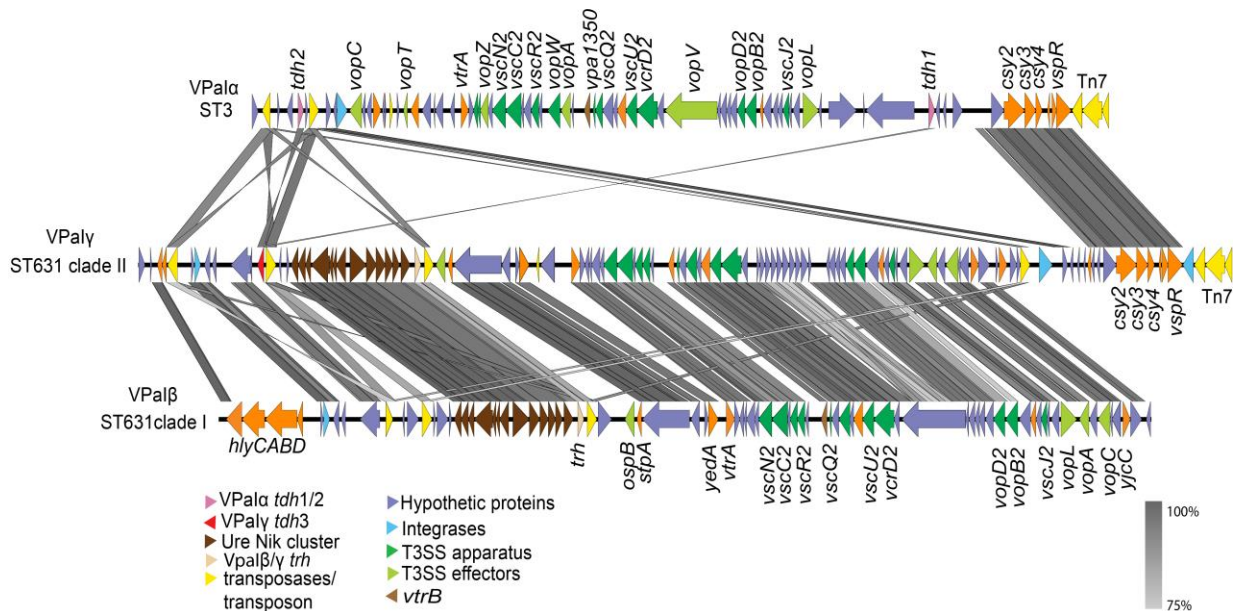


Figure 4. Comparisons of the pathogenicity islands containing hemolysins and T3SS2. Sequences of VP $\alpha$  were extracted from select genomes and aligned. VP $\alpha$  was derived from ST3 strain RIMD2210633, VP $\gamma$  was derived from ST631 clade II isolate MAVP-Q, and VP $\beta$  was derived from ST631 clade I isolate MAVP-R. ORFs are depicted in defined colors and similarities ( $\geq 75\%$ ) among ORFs are illustrated in grey blocks. omologs between VP $\alpha$  and VP $\beta/\gamma$  ( $50 > 75\%$  identity) are named and listed in supplemental table 2.

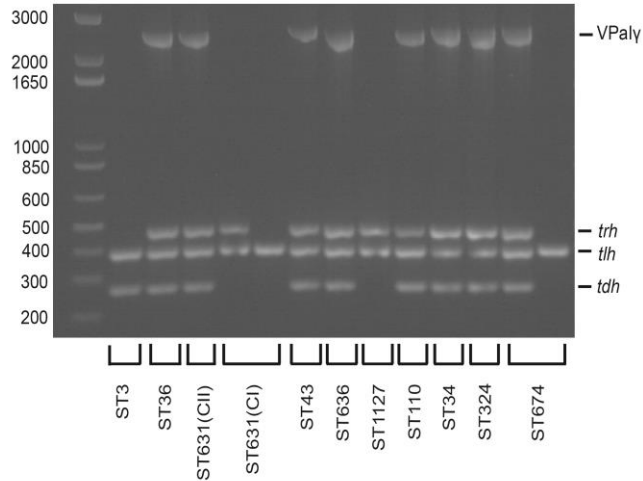


Figure 5. Distribution of VPaly in emergent pathogen lineages. The presence of *tdh*, *trh* and VPaly along with positive control *tlh* was determined by PCR amplification using gene-specific primers and visualized on a 1.2% agarose gel. The order from left to right is 1kb+ ladder, ST3 (MAVP-C), ST36 (MAVP-26), ST631 CII (clade II isolate MAVP-Q), ST631 CI (clade I isolates MAVP-R and G149), ST43 (MAVP-71), ST636 (MAVP-50), ST1127 (MAVP-M), ST110 (MAVP-46), ST34 (CTVP19C), ST324 (MAVP-14), and ST674 (CT4291, MAVP21). The corresponding sizes of the ladder fragments are as labeled to the left and the identity of the amplicons listed to the right of the gel image.