

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 lineages (ST36, ST43 and ST636) to  
26 near-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, a human  
29 pathogenic variant evolved independently through acquisition of distinct *Vibrio* pathogenicity  
30 islands (VPaI). These VPaI differ from each other and bear little resemblance to hemolysin-  
31 containing VPaI from isolates of the pandemic clonal complex. Clade I ST631 isolates either  
32 harbored no hemolysins, or contained a chromosome I-inserted island we call VPaI $\beta$  that  
33 encodes a type three secretion system (T3SS2 $\beta$ ) typical of Trh hemolysin-producers. The more  
34 clinically prevalent and clonal ST631 clade II had an island we call VPaI $\gamma$  that encodes both *tdh*  
35 and *trh* and that was inserted in chromosome II. VPaI $\gamma$  was derived from VPaI $\beta$  but with some  
36 additional acquired elements in common with VPaI carried by pandemic isolates, exemplifying  
37 the mosaic nature of pathogenicity islands. Genomics comparisons and amplicon assays  
38 identified VPaI $\gamma$ -type islands containing *tdh* inserted adjacent to the *ure* cluster in the three  
39 introduced Pacific and most other emergent lineages. that collectively cause 67% of Northeast  
40 US infections as of 2016.

41

42 **IMPORTANCE**

43           The availability of three different hemolysin genotypes in the ST631 lineage provided a  
44 unique opportunity to employ genome comparisons to further our understanding of the processes  
45 underlying pathogen evolution. The fact that two different pathogenic clades arose in parallel  
46 from the same potentially benign lineage by independent VP*a*I acquisition is surprising  
47 considering the historically low prevalence of community members harboring VP*a*I in waters  
48 along the Northeast US Coast that could serve as the source of this material. This illustrates a  
49 possible predisposition of some lineages to not only acquire foreign DNA but also to become  
50 human pathogens. Whereas the underlying cause for the expansion of *V. parahaemolyticus*  
51 lineages harboring VP*a*I $\gamma$  along the US Atlantic coast and spread of this element to multiple  
52 lineages that underlies disease emergence is not known, this work underscores the need to define  
53 the environment factors that favor bacteria harboring VP*a*I in locations of emergent disease.

54

## 55 **INTRODUCTION**

56           *Vibrio parahaemolyticus* is an emergent pathogen capable of causing human gastric  
57 infections when consumed, most often with contaminated shellfish (1, 2). Some human  
58 pathogenic *V. parahaemolyticus* variants evolve from diverse non-pathogenic communities  
59 through horizontal acquisition of *Vibrio* pathogenicity islands (VP*a*I) (3-5). Gastric pathogenic *V.*  
60 *parahaemolyticus* typically harbor islands with at least one of two types of horizontally acquired  
61 hemolysin genes (*tdh* and *trh*) that are routinely used for pathogen discrimination even though  
62 their role in disease appears modest (6-11). Most pathogenic *V. parahaemolyticus* isolates also  
63 carry accessory type three secretion systems (T3SS) that translocate effector proteins that  
64 contribute to host interaction (12-14). Two evolutionarily divergent horizontally-acquired  
65 accessory systems (T3SS2 $\alpha$  or T3SS2 $\beta$ ) contribute to human disease and are genetically linked

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

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

89 (the northwest Atlantic Ocean) (31). However, we previously identified a single ST631 isolate  
90 lacking hemolysins (21, 27) suggesting this pathogen lineage may have recently evolved through  
91 VP*aI* acquisition.

92 The goal of our study was to understand the genetic events and changing population  
93 context for the evolution of the ST631 pathogenic lineage. We conducted whole and core  
94 genome phylogenetic analysis of three environmental and 39 clinical ST631 isolates along with  
95 isolates from other emergent lineages from the region, which revealed two ST631 clades of  
96 common ancestry, from which human pathogens evolved in parallel. The single clade I clinical  
97 isolate acquired a *recA* gene insertion previously seen associated with Asian lineages, and had a  
98 VP*aI* that is typical of isolates harboring *trh* in the absence of *tdh*. In contrast, isolates from the  
99 clonal ST631 clade II that dominates Atlantic-derived ST631 infections (31) had a related but  
100 distinct VP*aI*. This VP*aI* contained a *tdh* gene and four associated hypothetical protein encoding  
101 genes inserted within, not next to, an existing *ure-trh-T3SS2 $\beta$*  island in close proximity to the  
102 *ure* cluster. Nearly all emergent resident and invasive lineages, including all three Pacific  
103 lineages (ST36, ST636 and ST43) contained islands that similarly had a *tdh* gene inserted within  
104 the VP*aI* in an identical location adjacent to the *ure* cluster providing a mechanism for  
105 simultaneous acquisition of both hemolysins with T3SS2 $\beta$ .

106

## 107 **RESULTS**

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

110 Ongoing analysis of clinical isolates revealed that even as the Pacific-derived ST36  
111 lineage continued to dominate infections (50%), the endemic (autochthonous) ST631 lineage

112 accounted for 14% of infections (Table 1). Concurrently, a limited number of other lineages  
113 contributed individually to fewer infections ( $\leq 3\%$  each), among which were two lineages that  
114 have caused infections in the Pacific Northwest in prior decades: ST43 and ST636 (22, 23).  
115 ST43 and ST636 only recently (2013 and 2011 respectively) (21) have been linked to product  
116 harvested from waters along the Northeast US coast, and also caused infections in subsequent  
117 years. As is common among US clinical isolates, pathogenic isolates of all the aforementioned  
118 lineages harbor both the *tdh* and *trh* hemolysin genes (Table 1). Among environmental isolates,  
119 ST34 and ST674 are the most frequently recovered pathogen lineages but these caused  
120 comparatively few infections (Table 1). ST34 was first reported from the environment in 1998,  
121 from both the Gulf of Mexico and near-shore areas of MA, and was also recovered in NH in  
122 2012 (21) suggesting it is an established resident in the region. ST674 which was first reported  
123 from an infection in Virginia in 2007 (32) was first recovered from the local environment in  
124 2012 ([www.pubmlst.org/vparahaemolyticus](http://www.pubmlst.org/vparahaemolyticus)) (21). Notably even though all four ST674  
125 environmental isolates, like ST34, harbored both hemolysin genes, the single ST674 clinical  
126 isolate (MAVP-21) lacked hemolysins (Table 1) (21). The decrease in clinical prevalence of *trh*-  
127 harboring Atlantic endemic ST1127, which caused no infections in the last three years, coincided  
128 with the increase in clinical prevalence of all three Pacific-derived lineages which harbor both  
129 hemolysins. Notably, very few other clinical isolates harbored *trh* in the absence of *tdh* and  
130 clinical isolates containing only *tdh* (i.e. ST1725) were extremely rare (Table 1). Concurrent with  
131 this shift in composition of clinical lineages that includes multiple Pacific-derived lineages,  
132 hemolysin producers have increased in relative abundance in nearshore areas of the region,  
133 where historically these represented  $\sim 1\%$  of all isolates (27). Since 2012, hemolysin producers  
134 have been recovered more frequently, and in the last two years their proportion has increased by

135 up to an order of magnitude (comprising as much as 10%) in some regional shellfish associated  
136 populations (data not shown).

137

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

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

152

153 **ST631 forms two divergent clades**

154       The existence of three different hemolysin profiles (Table 1) among all available ST631  
155 draft genomes suggested there could be more than one ST631 lineage. Therefore we evaluated  
156 whole genome maximum likelihood (ML) phylogenies of select ST631 isolates and all other  
157 lineages causing two or more infections reported in four Northeast US States to evaluate whether

158 there was more than one ST631 lineage (Table 1) (Fig. 2). The phylogenetic tree showed that  
159 ST631 isolates, regardless of their hemolysin genotype, clustered together but they formed two  
160 distinct clades, indicative of common ancestry (Fig. 2). Clade I harbored either *trh* or no  
161 hemolysins and consisted of all three environmental isolates which were from MA and NH, and  
162 the single clinical isolate MAVP-R, whereas clade II consisted of all other isolates all of which  
163 harbor both hemolysins. The two distinct ST631 clades shared 85% of their DNA in common  
164 and displayed polymorphisms in  $\leq 12\%$  of the shared DNA content. The most closely related  
165 sister lineage to ST631 was formed by *trh*-harboring ST1127 isolates that have been exclusively  
166 reported from clinical sources in the Northeast US (21).

167 We next evaluated the relationships of all available ST631 isolate genomes at NCBI and  
168 sequenced by us (Supplemental Table 1) using a custom core genome multi-locus sequence  
169 typing (cgMLST) method as previously described (31). Minimum spanning trees built from core  
170 genome loci from 42 ST631 isolates indicated that only 390 loci varied between the most closely  
171 related isolate of clade I (MAVP-L) and clade II (G6928) (Fig. 3). The most distantly related  
172 isolates within clade I (G149 and MAVP-R) exhibited 80 core genome loci differences whereas  
173 clade II is clonal with only 51 variant loci between the most divergent isolates: clinical isolate  
174 09-4436 and environmental isolate S487-4, both reported from PEI Canada (Fig. 3) (31).

175

176 **Each ST631 clade independently acquired a distinct pathogenicity island positioned on**  
177 **different chromosomes**

178 Given the variation in ST631, comparisons between these isolates could elucidate the  
179 events that led not only to the evolution of two pathogenic clades but also address unresolved  
180 questions about the unique configurations and contents of pathogenicity islands in western

181 Atlantic Ocean emergent lineages. The physical proximity of *tdh* with the *ure* cluster and *trh*,  
182 and the co-occurrence of *tdh* with T3SS2 $\beta$  reported in many *tdh*<sup>+</sup>/*trh*<sup>+</sup> clinical isolates suggested  
183 *tdh* could be harbored within or next to the same pathogenicity island harboring *trh* in at least  
184 some lineages as was previously suggested (20, 24, 34).

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

204 previously (16, 17, 20), VP*α*β is dissimilar to VP*α*α in most gene content with ~ 78 ORFs  
205 unique to VP*α*β (where the number of identified ORFs used for comparison can differ slightly  
206 depending on which annotation program is applied) (Supplemental Table 2, Fig. 4). Even so,  
207 VP*α*β had many homologous genes of varying sequence identity (n=~38 ORFs, excluding *tdh*  
208 homology with *trh*) when compared to VP*α*α (Supplemental Table 2, Fig. 4)(4, 5, 16).  
209 Identification of some homologs required that we relax matching to 50% such as for the  
210 divergent, but homologous T3SS2*α* and T3SS2*β* genes encoding the apparatus, chaperones, and  
211 some shared effectors (Supplemental Table 2). No homolog of the T3SS2*α* effector gene *vopZ*  
212 was identified, but a single ORF whose deduced protein sequence bears only 27% identity with  
213 VopZ is located in its place (Fig. 2 and Supplemental Table 2). VP*α*β from strain TH3996 and  
214 VP*α*α from pandemic strain RIMD 2210633 are inserted in an identical location in chromosome  
215 II adjacent to an Acyl-CoA hydrolase-encoding gene. In contrast the VP*α*βs in MAVP-R,  
216 ST1127 isolate MAVP-25, and Asia-derived AQ4037 are in chromosome I, in each case in the  
217 same insertion location identified for strain AQ4037 (17).

218       MAVP-Q contained both *tdh* and *trh* within the same contiguous unique VP*α*I (hereafter  
219 called VP*α*Iγ) that shared features with both VP*α*α and VP*α*β (Fig. 4, Supplemental Table 2).  
220 Specifically, VP*α*Iγ had a core that with few exceptions was orthologous in content and  
221 syntenous with VP*α*β from MAVP-R (Fig. 4) with only a few exceptions. VP*α*Iγ displays high  
222 conservation with VP*α*α near its 3' end, as has been described in other draft *tdh*<sup>+</sup>*trh*<sup>+</sup> harboring  
223 genomes (20) as well as in the VP*α*β island of strain TH3996, although the presence of this  
224 element may not be typical of VP*α*β (e.g. it is absent in the islands from AQ4037 (17), MAVP-  
225 R and MAVP-25). The VP*α*Iγ also contained a *tdh* gene homologous to *tdh2* (also called *tdhA*)  
226 from VP*α*α (98.6%) near its 5' end but not at the 5' terminus of the island (Fig. 4). Rather, the

227 DNA flanking both sides of the *tdh* gene in VPα $\gamma$  was conserved in VPα $\beta$  of MAVP-R and  
228 absent from VPα $\alpha$ , (Fig. 4). Analysis of 300 genomes of *V. parahaemolyticus* (representing a  
229 minimum of 28 distinct sequence types) of sufficient quality for analysis confirmed that the  
230 module of four hypothetical proteins preceding the *tdh2* homolog was present only in *trh*-  
231 harboring genomes, but not in genomes harboring *tdh* in the absence of *trh* (i.e. VPα $\alpha$   
232 containing genomes), providing evidence that the *tdh* gene was acquired horizontally by  
233 insertion into, not next to, an existing VPα $\beta$ , perhaps through activity of the adjacent transposase  
234 gene (11) (Supplemental Table 3, Supplemental fig. 1, and data not shown). Like with VPα $\alpha$   
235 from RIMD 2210633, and VPα $\beta$  of TH3996, VPα $\gamma$  of clade II ST631 is located in a conserved  
236 location of chromosome II, adjacent to an Acyl-CoA hydrolase-encoding gene.

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

243

244 **Most clinically prevalent isolates from the Northeast US harbor similar contiguous**  
245 **pathogenicity islands containing *tdh* inserted in the same location of their VPαI**

246         We next asked which isolates from other lineages likely residing within the mixed  
247 population with ST631 in near-shore areas of the Northeast US harbored islands of similar  
248 structure to VPα $\gamma$  that contain both hemolysin genes. Assembly of short-read sequences into  
249 contigs that cover the full length of VPαI which is necessary for comparative analysis of entire

250 island configuration was impeded by the fact that homologous transposase sequences and other  
251 sequences were repeated multiple times throughout the island. Therefore, we determine whether  
252 other lineages harboring both hemolysin genes harbor *tdh* in the same island location, between  
253 the conserved VP $\alpha$ I $\beta$ / $\gamma$  module of four hypothetical proteins (to the left or 5' of *tdh*) and the *ure*  
254 cluster (to the right or 3' of *tdh*) (Fig. 4) by combining bioinformatics analysis of sequenced  
255 genomes with amplicon assays (Supplemental Fig. 1). First we analyzed assembled draft  
256 genomes for *tdh* co-occurrence and proximity with the four adjacent hypothetical protein-  
257 encoding genes that are absent in VP $\alpha$ I $\beta$  but present in VP $\alpha$ I $\gamma$  (See Methods). Every emergent  
258 pathogenic lineage of the Northeast US (Table 1) harboring both *tdh* and *trh* carried homologous  
259 DNA corresponding to all four hypothetical proteins adjacent to the *tdh* gene in a contiguous  
260 segment (Supplemental Table 3). To determine whether *tdh* was also adjacent to the *ure* cluster  
261 in these same isolates we next designed specific flanking primers and amplified the unique  
262 juncture between the *tdh*-containing transposon associated module and the *ure* cluster for all  
263 clinical isolates harboring both *tdh* and *trh* (See Methods) (Supplemental Fig. 1). The results  
264 were congruent with our bioinformatics assessment (Supplemental Table 3), and demonstrated  
265 that isolates from all emergent pathogenic lineages harboring both hemolysins have *tdh* inserted  
266 in close proximity to an *ure* cluster in a configuration similar to VP $\alpha$ I $\gamma$  from MAVP-Q (Fig. 5,  
267 Table 1). This confirmed that these isolates harboring both hemolysins harbor *tdh* within, and not  
268 next to, the same VP $\alpha$ I thereby facilitating simultaneous acquisition of both hemolysin genes.

269

## 270 **DISCUSSION**

271 Even preceding the increased illnesses from Pacific-invasive lineages, two different  
272 clades of the predominant endemic Atlantic lineage of pathogenic *V. parahaemolyticus*, ST631

273 (31) evolved and contributed to a rise in sporadic illnesses in the four reporting Northeast US  
274 States (Table 1, Fig. 2 & 3). Several lines of evidence support the interpretation of parallel  
275 pathogen evolution. The two lineages exhibit differences in both clinical and environmental  
276 prevalence suggesting the pathogenic variants of each clade have not evolved the same degree of  
277 virulence (Table 1). Pathogenic members in each lineage also acquired different pathogenicity  
278 islands with different hemolysin gene content (Fig. 2 & 3). Although it was a formal possibility  
279 that ST631 clade II evolved from clade I by independent horizontal acquisition of *tdh* into its  
280 existing VPαIβ, it is notable that other resident and even invasive lineages now in the Atlantic  
281 harbor VPαIγ with *tdh* and four additional co-occurring ORFs inserted into the same location of  
282 the island, suggesting a common evolutionary origin of this hybrid-type island (Fig. 4 and  
283 Supplemental Fig. 1). Finally, each of the two clades harbor VPαI insertions on different  
284 chromosomes: the less clinically prevalent ST631 clade I contains three isolates that harbor  
285 VPαIβ in chromosome I (Fig. 3) and a single environmental isolate lacking any island (Table 1,  
286 supplemental Fig. 2), whereas the clonal ST631 clade II isolates all harbor VPαIγ on  
287 chromosome II.

288         Given that several other resident lineages harbor similar β and γ-type VPαI, pathogens in  
289 each clade could have acquired their islands from the reservoir of resident bacteria already  
290 circulating in the Atlantic even before the presume arrival of invasive Pacific lineages. Several  
291 well-documented members of the Gulf of Mexico *V. parahaemolyticus* population (35-37) may  
292 also have expanded their range through movement of ocean currents and could be the source for  
293 these VPαI (Table 1, Fig. 5). But historically, hemolysin producers were extremely rare in near  
294 shore areas of the Atlantic US coast (25) and represented only about ~1% of isolates in an  
295 estuary of NH as of a decade ago (27) limiting the potential for interacting partners or sources for

296 acquired VP*a*I. Given this historical context, it is remarkable that two different clades from the  
297 same lineage independently acquired different VP*a*I-which for clade II ST631 occurred prior to  
298 2007 -well before the recent shift in abundance of hemolysin producers.

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

318 lineages are warranted given this evidence and the documented introduction of multiple Pacific-  
319 derived lineages in the region (Table 1).

320         There has been some consideration of the roles of human virulence determinants in  
321 ecological fitness, but the natural context of pathogenic *V. parahaemolyticus* evolution is still  
322 unknown (40-42). Whereas *tdh* and T3SS2 $\alpha$  each may promote growth when bacteria are under  
323 predation, isolates that carry *trh*-containing islands (which likely also have T3SS2 $\beta$ ) do not  
324 derive similar benefits from their islands (43). This is surprising considering the islands encode  
325 several homologous effectors (Fig. 4 and Supplemental Table 2) that don't have an established  
326 role in enteric disease but they could alternatively or additionally mediate eukaryotic cell  
327 interactions with natural hosts thereby promoting environmental fitness (13, 14). But these  
328 islands also lack homologous open for the VP $\alpha$  effector that is most closely associated with  
329 enteric disease: *vopZ* (11) (Fig. 4 and Supplemental Table 2). The general lack of knowledge of  
330 unique T3SS2 $\beta$  effectors and other gene function in these islands (Fig. 4 and Supplemental Table  
331 2) even with regard to enteric disease, limits comparative analysis with the well-studied and  
332 functionally defined VP $\alpha$  which could elucidate the bases for pathogen evolution. The higher  
333 clinical prevalence of clade II ST631 than clade I which has also been recovered on more than  
334 one occasion from the environment (Table 1) could indicate that VP $\gamma$  confers greater virulence  
335 potential than VP $\beta$ , perhaps owing to the presence of *tdh*, a known virulence factor (1, 7, 44).  
336 However, the resident community members in both the Pacific and the Atlantic Ocean that  
337 harbor *tdh* and T3SS2 $\alpha$  comparatively rarely cause human infections (21-23). The unique  
338 environmental conditions that underlie pathogen success from northern latitudes that favors  
339 bacteria with VP $\beta$  and VP $\gamma$  including two different ST631 lineages suggests the shared  
340 content of these islands could confer abilities that are distinct from VP $\alpha$  which could underlie

341 the repeated acquisition and maintenance of these related islands by so many different lineages  
342 now present in near-shore areas of the Northeast US.

343

## 344 **MATERIALS AND METHODS**

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

346 *V. parahaemolyticus* clinical isolates for this study were provided by cooperating public  
347 health laboratories in Massachusetts, New Hampshire, Maine, and Connecticut whereas a select  
348 number of environmental isolates were enriched from estuarine substrates as described (21).  
349 Detailed information about these isolates was described previously (31) and listed in  
350 Supplemental Table 1. Isolates were routinely cultured in Heart Infusion (HI) media  
351 supplemented with NaCl at 37°C as described (21).

352

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

354 Genomic DNA was extracted using the Wizard Genomic DNA purification Kit (Promega,  
355 Madison WI USA) or by organic extraction (21). The quality genomic DNA was determined by  
356 spectrophotometric measurements by NanoDrop (ThermalFisher, Waltham MA USA). Libraries  
357 for DNA sequencing were prepared using a high-throughput Nextera DNA preparation protocol  
358 (45) using an optimal DNA concentration of 2ng/μl. Genomic DNA was sequenced using an  
359 Illumina – HiSeq2500 device at the Hubbard Center for Genome Studies at the University of  
360 New Hampshire, using a 150bp paired-end library. *De novo* assembly was performed using the  
361 A5 pipeline (46), and the assemblies annotated with Prokka1.9 using the "genus" option and  
362 selecting "*Vibrio*" for the reference database (47). The sequence types were subsequently  
363 determined using the SRST2 pipeline (48). The sequence type of each genome was determined

364 when using *V. parahaemolyticus* as the database (<https://pubmlst.org/vparahaemolyticus/>). For  
365 most isolates where the combination of each allele was not found in the database representing  
366 novel sequence types, the genome was submitted for a new sequence type designation  
367 ([www.pubmlst.org/vparahaemolyticus](http://www.pubmlst.org/vparahaemolyticus)).

368 Isolates MAVP-Q and MAVP-R were sequenced using the Pacific Biosciences RSII  
369 technology. Using between 3.7-5.3  $\mu\text{g}$  DNA, the library preparation and sequencing was  
370 performed according to the manufacturer's instructions (Pacific Biosciences, Menlo Park CA,  
371 USA) and reflects the P5-C3 sequencing enzyme and chemistry for MAVP-Q isolate and the P6-  
372 C4 configuration for MAVP-R. The mass of double-stranded DNA was determined by Qubit  
373 (Waltham, MA USA) and the sample diluted to a final concentration of 33  $\mu\text{g} / \mu\text{L}$  in a volume  
374 of 150  $\mu\text{L}$  elution buffer (Qiagen, Germantown MD USA). The DNA was sheared for 60  
375 seconds at 4500 rpm in a G-tube spin column (Covaris, Woburn MA USA) which was  
376 subsequently flipped and re-spun for another 60 seconds at 4500 rpm resulting in a ~20,000 bp  
377 DNA verified using a DNA 12000 Bioanalyzer gel chip (Agilent, Santa Clara, CA USA). The  
378 sheared DNA isolate was then re-purified using a 0.45X AMPure XP purification step (Beckman  
379 Coulter, Indianapolis IN USA). The DNA was repaired by incubation in DNA Damage Repair  
380 solution. The library was again purified using 0.45X Ampure XP and SMRTbell adapters ligated  
381 to the ends of the DNA at 25°C overnight. The library was treated with an exonuclease cocktail  
382 (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  
383 fragments. Two additional 0.45X Ampure XP purifications steps were performed to remove  
384 <2000 bp molecular weight DNA and organic contaminant.

385 Upon completion of library construction, samples were validated using an Agilent  
386 DNA 12000 gel chip. The isolate library was subjected to additional size selection to the range

387 of 7,000 bp – 50,000 bp to remove any SMRTbells < 5,000 bp using Sage Science Blue Pippin  
388 0.75% agarose cassettes to maximize the SMRTbell sub-read length for optimal *de*  
389 *novo* assembly. Size-selection was confirmed by Bio-Analysis and the mass was quantified using  
390 the Qubit assay. Primer was then annealed to the library (80°C for 2 minute 30 followed by  
391 decreasing the temperature by 0.1°/s to 25°C). The polymerase-template complex was then  
392 bound to the P5 or P6 enzyme using a ratio of 10:1 polymerase to SMRTbell at 0.5 nM for 4  
393 hours at 30°C and then held at 4°C until ready for magbead loading, prior to sequencing. The  
394 magnetic bead-loading step was conducted at 4°C for 60-minutes per manufacturer’s guidelines.  
395 The magbead-loaded, polymerase-bound, SMRTbell libraries were placed onto the RSII machine  
396 at a sequencing concentration of 110-150 pM and configured for a 180-minute continuous  
397 sequencing run. Long read assemblies were constructed using HGAP version 2.3.0 for *de novo*  
398 assembly generation. Further, hybrid assemblies were generated and error corrected with  
399 illumina raw reads using Pilon v1.20 (49).

400

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

402 To more rapidly identify ST631 isolates from clinical and environmental collections we  
403 developed PCR-amplicon assays to unique gene content in ST631. Whole genome comparisons  
404 were performed on MAVP-Q (a ST631 clinical isolate), G149 (a ST631 environmental isolate),  
405 MAVP-26 (ST36), RIMD2210633 (ST3), and AQ4037 (ST96) (Supplemental Fig. 3). A total of  
406 26 distinct genomic regions, each greater than 1kb in size, were present in MAVP-Q but absent  
407 in other comparator genomes, including environmental ST631 that lacks hemolysins (G149)  
408 (Supplemental Fig. 3). Within a large genomic island ~37.6 Kb in length with an integrase at one  
409 terminus and an overall lower GC content (40.6% compared to 45.8% for the genome) a single

410 ORF homologous to restriction endonucleases (AB831\_06355) that was restricted to clinical  
411 ST631 isolates in our collection and publicly available draft genomes (n=693)  
412 (<http://www.ncbi.nlm.nih.gov/genome/691>, 2017) was selected as a suitable amplicon target. The  
413 distribution of this locus was further analyzed using the BLAST algorithm by a query against the  
414 nucleotide collection, the non-redundant protein sequences, and against the genus *Vibrio* (taxid:  
415 662), excluding *V. parahaemolyticus* (taxid: 691), using the default settings for BLASTn (50).  
416 Similar approaches were applied to identify ST631 diagnostic loci inclusive of the single  
417 environmental isolate (G149), which identified a hypothetical protein encoding region  
418 (AB831\_06535) (ST631env). Oligonucleotide primers were designed to amplify the diagnostic  
419 regions including AB831\_06355 using primers ST631end F  
420 (5'AGTTCATCAGGTAGAGAGTTAGAGGA3') and ST631endR  
421 (5'TCTTCGTTACCATAGTATGAGCCA3') which produces an amplicon of c.a. 494bp, and  
422 AB831\_06535 using primers ST631envF (5'TGGGCGTTAGGCTTTGTC3') and ST631-envR  
423 (5'GGGCTTCTACGACTTTCTGCT3') producing an amplicon of 497bp.

424       Amplification of diagnostic loci was evaluated in individual assays using genomic DNA  
425 from positive and negative controls: MAVP-Q and G149 (ST631), G4186 (ST34), G3578  
426 (ST674), and MAVP-M (ST1127), MAVP-26 (ST36) and G61 (ST1125). Amplification of  
427 specific sequence types were performed with Accustart enzyme mix on purified DNA. Cycling  
428 was performed with an initial denaturation at 94°C for 3 min., followed by 30 cycles of a  
429 denaturation at 94°C for 1min, annealing at 55°C for 1 min, and amplification at 72°C for 30s  
430 with a final elongation at 72°C for 5 min. The primer pairs only produced amplicons from  
431 template DNA from ST631 and each was the expected size (data not shown, and Supplemental  
432 Fig. 3). Amplicon assays were applied to 208 clinical isolates from the Northeast US States (ME,

433 NH, MA and CT) and 1140 environmental isolates collected from 2015-2016 from NH and MA.  
434 These assays identified all known ST631 clinical isolates with 100% specificity and also  
435 identified an additional 7 *tdh*<sup>+</sup>*trh*<sup>+</sup> clinical isolates (ST631*end* and ST631*env* positive), and two  
436 environmental (ST631*end* negative and ST631*env* positive) isolates from our archived collection.  
437 Each, with the exception of MAVP-R, was subsequently confirmed to be ST631 by seven-locus  
438 MLST ([www.pubmlst.org](http://www.pubmlst.org)).

439

#### 440 **Examination of *recA* allele and adjacent sequences**

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

447

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

449 Whole genome phylogenies were constructed with single nucleotide polymorphisms  
450 (SNPs) identified from draft genomes using kSNP3 to produce aligned SNPs in FASTA format  
451 (52). A maximum likelihood (ML) tree was then built from the FASTA file using raxMLHPC  
452 with model GTRGAMMA and the -f option, and 100 bootstraps (53). Since there were no  
453 differences among the clade II ST631 isolates we used a subset representing geographic and  
454 temporal span of isolation.

455 Minimum spanning tree (MST) analysis was built based on core gene SNPs produced  
456 from a cluster analysis. The cluster analysis of ST631 was performed using a custom core  
457 genome multi-locus sequence type (cgMLST) analysis using RidomSeqSphere+software v3.2.1  
458 (<http://www.ridom.de.seqsphere>, Ridom GmbH, Münster, Germany) as previously described  
459 (31). Briefly, the software first defines a cgMLST scheme using the target definer tool with  
460 default settings using the PacBio generated MAVP-Q genome as the reference. Then, five other  
461 *V. parahaemolyticus* genomes (BB22OP, CDC\_K4557, FDA\_R31, RIMD2210633, and UCM-  
462 V493) were used for comparison with the reference genome to establish the core and accessory  
463 genome genes. Genes that are repeated in more than one copy in any of the six genomes were  
464 removed from the analysis. Subsequently, a task template was created that contains both core and  
465 accessory genes. Each individual gene locus from MAVP-Q was assigned allele number 1. Then  
466 each ST631 isolate genome assembly was queried against the task template, where any locus that  
467 differed from the reference genome or any other queried genome was assigned a new allele  
468 number. The cgMLST performed a gene-by-gene analysis of all core genes (excluding accessory  
469 genes) and identified SNPs within different alleles to establish genetic distance calculations.

470

#### 471 **Configuration and distribution of VPais**

472 The VPai sequence from the PacBio sequenced genomes of MAVP-Q and MAVP-R  
473 were identified by comparison with the published RIMD2210633 VPai-7 (NC\_004605 region  
474 between VPA1312 – VPA1395) and VPai<sub>TH396</sub> (AB455531) (16). Identification of the complete  
475 MAVP-Q VPai<sub>γ</sub> and genomic junctures in chromosome II was done by comparison with the  
476 same region of chromosome II in MAVP-R and G149 (which lack an island in this location)  
477 using Mauve (54). In a reciprocal manner, the absence of an island in chromosome I in MAVP-Q

478 and G149 was assessed by comparison with chromosome I of MAVP-R. MAVP-Q VP*AI* $\gamma$   
479 (MF066646) and MAVP-R VP*AI* $\beta$  (MF066647) were then extracted as a single contiguous  
480 sequence and annotated using Prokka 1.9. Gene content and order of the VP*AI* elements in  
481 MAVP-Q, MAVP-R and RIMD2210633 were then illustrated by Easyfig (51). Roary (55) was  
482 then employed to determine homologs among VP*AI*s based on each island's annotated sequences  
483 with identity set at 50%. Identification of the genome locations of VP*AI* $\beta$  in ST1127 isolate  
484 MAVP-M (accession number GCA\_001023155) and for VP*AI* $\gamma$  in AQ4037 (accession number  
485 GCA\_000182365) (17) was also done using Mauve (54).

486 To examine the distribution of the VP*AI* $\gamma$  in all publicly available draft genomes  
487 (<https://www.ncbi.nlm.nih.gov/genome/genomes/691>, 2016) and genomes from archived  
488 regional isolates, whole draft genome sequences were aligned to a 6,118 bp subsequence of the  
489 MAVP-Q VP*AI* with NASP version 1.0.2 (56) (<https://pypi.python.org/pypi/nasp/1.0.2>, 2017).  
490 This subsequence spanned the unique juncture of the four conserved hypothetical proteins  
491 (AB831\_22090, AB831\_22095, AB831\_22100, AB831\_22105) with the adjacent inserted *tdh*  
492 (AB831\_22110, c.a. 2549 bp upstream of *ure* cluster)(Supplemental Fig. 1). Percent coverage of  
493 the reference sequence was used to determine whether each genome harbored only the four  
494 hypothetical proteins, only a *tdh* gene, or the entire module including the fusion of the four genes  
495 with *tdh* (Supplemental Fig. 1 and Supplemental Table 3). The sequence type of each genome  
496 harboring the fused element characteristic of VP*AI* $\gamma$  was then determined using the SRST2  
497 pipeline (48). Where sequencing reads were not available as the input for SRST2, they were  
498 simulated from assemblies using an in-house Python script  
499 (<https://github.com/kpdrees/fasta2reads>).

500 A PCR amplification approach was developed and applied to survey the presence of *tdh*  
501 adjacent to the *ure* gene cluster. Primers were designed to conserved sequences of the 3' end of  
502 *tdh* (PIHybF8: 5'GCCAACATGGATATAAATAAAAATGA3') and the 5' end of *ureG*  
503 (*tdhUreGrev5*: 5'GACAAAGGTATGCTGCCAAAAGTG3') as determined by gene alignments,  
504 which when used together produced a 2631 bp amplicon of the insertion juncture when used with  
505 MAVP-Q as a template (Supplemental Fig. 4). Amplification was performed on purified DNA  
506 with Accustart enzyme mix, with an initial denaturation at 94°C for 3 min., followed by 30  
507 cycles of a denaturation at 94°C for 1 min, annealing at 61°C for 1min, and amplification at 72°C  
508 for 2.5 min, with a final elongation at 72°C for 5 min. This amplification was performed in  
509 parallel with a diagnostic multiplex PCR amplification of *tdh*, *trh* and *tlh* using published  
510 methods (10, 57) to investigate the co-occurrence of VP*α**γ* with both hemolysin encoding genes  
511 in representative isolates of various clinically prevalent sequence types. Amplicons were  
512 visualized using a 1.2% agarose gel in TAE buffer (Supplemental Fig. 4).

513

#### 514 **Nucleotide sequence accession numbers.**

515 The accession number of Pacific Biosciences sequenced genome for MAVP-Q is  
516 MDWT000000000, and for MAVP-R is MPPP000000000. The accession number of Illumina  
517 sequenced draft genome for G6928 is MPPN000000000, for MA561 is MPPM000000000 and for  
518 G149 is MPPO000000000. Detailed information about all other ST631 isolate draft genomes were  
519 described previously (31) and are listed in Supplemental Table 1. The accessions for the short  
520 reads for the remaining sequenced genomes are listed in Supplemental Table 4. The accession  
521 number of VP*α**β* from MAVP-R is MF066647 and the accession number of VP*α**γ* from MAVP-  
522 Q is MF066646.

523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545

## ACKNOWLEDGEMENTS

We are grateful for clinical isolates and wish to thank specifically: Jana Ferguson and Tracy Stiles of the Massachusetts Department of Public Health, and M. Hickey and C. Schillaci from the Massachusetts Department of Marine Fisheries; J.K. Kanwit of the Maine Department of Marine Resources and A. Robbins from the Maine Department of Health and Human Services; and Larn Mank from the Connecticut Department of Public Health Laboratory, and K. DeRosia-Banick, Connecticut Department of Agriculture, Bureau of Aquaculture. Assistance with genome sequencing was provided by W. K. Thomas, and technical assistance provided by J. Lemaire, K. Hartman, C. Hallee, M. Malanga, S. Ilyas, J. Hall, J. Sevigny, M. Dillon, K. Flynn, A. Goupil, J. Means, R. Foxall, E. DaSilva, and M.S. Pankey. Partial funding for this work was provided by the USDA National Institute of Food and Agriculture (Hatch projects NH00574, NH00609 [accession number 233555], and NH00625 [accession number 1004199]). Additional funding was provided by the National Oceanic and Atmospheric Administration College Sea Grant program and grants R/CE-137, R/SSS-2, and R/HCE-3. Support was also provided through the National Institutes of Health (1R03AI081102-01), the National Science Foundation (EPSCoR IIA-1330641), and the National Science Foundation (DBI 1229361 NSF MRI). N.G.-E. was funded through the FDA Foods Science and Research Intramural Program. Feng Xu and Cheryl A. Whistler declare a potential conflict of interest in the form of a pending patent application (U.S. patent application 62/128,764). This is Scientific Contribution Number 2722 for the New Hampshire Agricultural Experiment Station.

546 **REFERENCES**

- 547 1. **Hiyoshi H, Kodama T, Iida T, Honda T.** 2010. Contribution of *Vibrio*  
548 *parahaemolyticus* virulence factors to cytotoxicity, enterotoxicity, and lethality in mice.  
549 *Infect Immun* **78**:1772-1780.
- 550 2. **Scallan E, Hoekstra RM, Angulo FJ, Tauxe RV, Widdowson M-A, Roy SL, Jones**  
551 **JL, Griffin PM.** 2011. Foodborne illness acquired in the United States—major  
552 pathogens. *Emerg Infect Dis* **17**(1):7-15.
- 553 3. **Hazen TH, Pan L, Gu J-D, Sobecky PA.** 2010. The contribution of mobile genetic  
554 elements to the evolution and ecology of *Vibrios*. *FEMS Microbiol Ecol* **74**:485-499.
- 555 4. **Hurley CC, Quirke A, Reen FJ, Boyd EF.** 2006. Four genomic islands that mark post-  
556 1995 pandemic *Vibrio parahaemolyticus* isolates. *BMC Genomics* **7**:104  
557 DOI:110.1186/1471-2164-1187-1104.
- 558 5. **Boyd EF, Cohen AL, Naughton LM, Ussery DW, Binnewies TT, Stine OC, Parent**  
559 **MA.** 2008. Molecular analysis of the emergence of pandemic *Vibrio parahaemolyticus*.  
560 *BMC Microbiol* **8**:110.
- 561 6. **Kishishita M, Matsuoka N, Kumagai K, Yamasaki S, Takeda Y, Nishibuchi M.** 1992.  
562 Sequence variation in the thermostable direct hemolysin-related hemolysin (*trh*) gene of  
563 *Vibrio parahaemolyticus*. *Appl Environ Microbiol* **58**:2449-2457.
- 564 7. **Honda T, Ni Y, Miwatani T, Adachi T, Kim J.** 1992. The thermostable direct  
565 hemolysin of *Vibrio parahaemolyticus* is a pore-forming toxin. *Can J Microbiol* **38**:1175-  
566 1180.

- 567 8. **Park K-S, Ono T, Rokuda M, Jang M-H, Iida T, Honda T.** 2004. Cytotoxicity and  
568 enterotoxicity of the thermostable direct hemolysin-deletion mutants of *Vibrio*  
569 *parahaemolyticus*. *Microbiol Immunol* **48**:313-318.
- 570 9. **Shirai H, Ito H, Hirayama T, Nakamoto Y, Nakabayashi N, Kumagai K, Takeda Y,**  
571 **Nishibuchi M.** 1990. Molecular epidemiologic evidence for association of thermostable  
572 direct hemolysin (TDH) and TDH-related hemolysin of *Vibrio parahaemolyticus* with  
573 gastroenteritis. *Infect Immun* **58**:3568-3573.
- 574 10. **Panicker G, Call DR, Krug MJ, Bej AK.** 2004. Detection of pathogenic *Vibrio* spp. in  
575 shellfish by using multiplex PCR and DNA microarrays. *Appl Environ Microbiol*  
576 **70**:7436-7444.
- 577 11. **Nishibuchi M, Kaper JB.** 1995. Thermostable direct hemolysin gene of *Vibrio*  
578 *parahaemolyticus*: a virulence gene acquired by a marine bacterium. *Infect Immun*  
579 **63**:2093.
- 580 12. **Park K-S, Ono T, Rokuda M, Jang M-H, Okada K, Iida T, Honda T.** 2004.  
581 Functional characterization of two type III secretion systems of *Vibrio parahaemolyticus*.  
582 *Infect Immun* **72**:6659-6665.
- 583 13. **Broberg CA, Calder TJ, Orth K.** 2011. *Vibrio parahaemolyticus* cell biology and  
584 pathogenicity determinants. *Microb Infect* **13**:992-1001.
- 585 14. **Zhang L, Orth K.** 2013. Virulence determinants for *Vibrio parahaemolyticus* infection.  
586 *Curr Opin Microbiol* **16**:70-77.
- 587 15. **Makino K, Oshima K, Kurokawa K, Yokoyama K, Uda T, Tagomori K, Iijima Y,**  
588 **Najima M, Nakano M, Yamashita A.** 2003. Genome sequence of *Vibrio*

- 589            *parahaemolyticus*: a pathogenic mechanism distinct from that of *V. cholerae*. The Lancet  
590            **361**:743-749.
- 591    16.    **Okada N, Iida T, Park K-S, Goto N, Yasunaga T, Hiyoshi H, Matsuda S, Kodama T,**  
592            **Honda T.** 2009. Identification and characterization of a novel type III secretion system in  
593            trh-positive *Vibrio parahaemolyticus* strain TH3996 reveal genetic lineage and diversity  
594            of pathogenic machinery beyond the species level. Infect Immun **77**:904-913.
- 595    17.    **Chen Y, Stine OC, Badger JH, Gil AI, Nair GB, Nishibuchi M, Fouts DE.** 2011.  
596            Comparative genomic analysis of *Vibrio parahaemolyticus*: serotype conversion and  
597            virulence. BMC Genomics **12**:1.
- 598    18.    **Zhou X, Gewurz BE, Ritchie JM, Takasaki K, Greenfeld H, Kieff E, Davis BM,**  
599            **Waldor MK.** 2013. *vopZ* A *Vibrio parahaemolyticus* T3SS effector mediates  
600            pathogenesis by independently enabling intestinal colonization and inhibiting TAK1  
601            activation. Cell Reports **3**:1690-1702.
- 602    19.    **Hubbard TP, Chao MC, Abel S, Blondel CJ, zur Wiesch PA, Zhou X, Davis BM,**  
603            **Waldor MK.** 2016. Genetic analysis of *Vibrio parahaemolyticus* intestinal colonization.  
604            Proc Nat Acad Sci USA **113**:6283-6288.
- 605    20.    **Ronholm J, Petronella N, Leung CC, Pightling A, Banerjee S.** 2016. Genomic  
606            Features of Environmental and Clinical *Vibrio parahaemolyticus* Isolates Lacking  
607            Recognized Virulence Factors Are Dissimilar. Appl Environ Microbiol **82**:1102-1113.
- 608    21.    **Xu F, Ilyas S, Hall JA, Jones SH, Cooper VS, Whistler CA.** 2015. Genetic  
609            characterization of clinical and environmental *Vibrio parahaemolyticus* from the  
610            Northeast USA reveals emerging resident and non-indigenous pathogen lineages. Name:  
611            Front Microbiol **6**:272.

- 612 22. **Banerjee SK, Kearney AK, Nadon CA, Peterson C-L, Tyler K, Bakouche L, Clark**  
613 **CG, Hoang L, Gilmour MW, Farber JM.** 2014. Phenotypic and genotypic  
614 characterization of Canadian clinical isolates of *Vibrio parahaemolyticus* collected from  
615 2000 to 2009. J Clin Microbiol **52**:1081-1088.
- 616 23. **Turner JW, Paranjpye RN, Landis ED, Biryukov SV, González-Escalona N, Nilsson**  
617 **WB, Strom MS.** 2013. Population structure of clinical and environmental *Vibrio*  
618 *parahaemolyticus* from the Pacific Northwest coast of the United States. PLoS ONE  
619 **8(2):e55726**
- 620 24. **Jones JL, Lüdeke CH, Bowers JC, Garrett N, Fischer M, Parsons MB, Bopp CA,**  
621 **DePaola A.** 2012. Biochemical, serological, and virulence characterization of clinical and  
622 oyster *Vibrio parahaemolyticus* isolates. J Clin Microbiol **50(7)**:2343-2352.
- 623 25. **DePaola A, Ulaszek J, Kaysner CA, Tenge BJ, Nordstrom JL, Wells J, Puhr N,**  
624 **Gendel SM.** 2003. Molecular, serological, and virulence characteristics of *Vibrio*  
625 *parahaemolyticus* isolated from environmental, food, and clinical sources in North  
626 America and Asia. Appl Environ Microbiol **69**:3999-4005.
- 627 26. **Haendiges J, Timme R, Allard MW, Myers RA, Brown EW, Gonzalez-Escalona N.**  
628 2015. Characterization of *Vibrio parahaemolyticus* clinical strains from Maryland (2012–  
629 2013) and comparisons to a locally and globally diverse *V. parahaemolyticus* strains by  
630 whole-genome sequence analysis. Front Microbiol **6**:125
- 631 27. **Ellis CN, Schuster BM, Striplin MJ, Jones SH, Whistler CA, Cooper VS.** 2012.  
632 Influence of seasonality on the genetic diversity of *Vibrio parahaemolyticus* in New  
633 Hampshire shellfish waters as determined by multilocus sequence analysis. Appl Environ  
634 Microbiol **78**:3778-3782.

- 635 28. **Nair GB, Ramamurthy T, Bhattacharya SK, Dutta B, Takeda Y, Sack DA.** 2007.  
636 Global dissemination of *Vibrio parahaemolyticus* serotype O3: K6 and its serovariants.  
637 Clin Microbiol Rev **20**:39-48.
- 638 29. **Martinez-Urtaza J, Baker-Austin C, Jones JL, Newton AE, Gonzalez-Aviles GD,**  
639 **DePaola A.** 2013. Spread of Pacific Northwest *Vibrio parahaemolyticus* strain. N Engl J  
640 Med **369**:1573-1574.
- 641 30. **Newton AE, Garrett N, Stroika SG, Halpin JL, Turnsek M, Mody RK, Division of**  
642 **Foodborne W, Environmental D.** 2014. Notes from the field: Increase in *Vibrio*  
643 *parahaemolyticus* infections associated with consumption of Atlantic coast shellfish—  
644 2013. MMWR Morb Mortal Wkly Rep **63**:335-336.
- 645 31. **Xu F, Gonzalez-Escalona N, Haendiges J, Myers RA, Ferguson J, Stiles T, Hickey E,**  
646 **Moore M, Hickey JM, Schillaci C.** 2017. Sequence type 631 *Vibrio parahaemolyticus*,  
647 an emerging foodborne pathogen in North America. J Clin Microbiol **55**:645-648.
- 648 32. **Lüdeke CH, Gonzalez-Escalona N, Fischer M, Jones JL.** 2015. Examination of  
649 clinical and environmental *Vibrio parahaemolyticus* isolates by multi-locus sequence  
650 typing (MLST) and multiple-locus variable-number tandem-repeat analysis (MLVA).  
651 Frontiers in microbiology **6**:564
- 652 33. **González-Escalona N, Gavilan RG, Brown EW, Martínez-Urtaza J.** 2015.  
653 Transoceanic spreading of pathogenic strains of *Vibrio parahaemolyticus* with distinctive  
654 genetic signatures in the recA gene. PloS one **10**:e0117485.
- 655 34. **Park K-S, Suthienkul O, Kozawa J, Yamaichi Y, Yamamoto K, Honda T.** 1998.  
656 Close proximity of the *tdh*, *trh* and *ure* genes on the chromosome of *Vibrio*  
657 *parahaemolyticus*. Microbiology **144**:2517-2523.

- 658 35. **Johnson C, Flowers A, Young V, Gonzalez-Escalona N, DePaola A, Noriega III N,**  
659 **Grimes D.** 2009. Genetic relatedness among *tdh+* and *trh+* *Vibrio parahaemolyticus*  
660 cultured from Gulf of Mexico oysters (*Crassostrea virginica*) and surrounding water and  
661 sediment. *Microb Ecol* **57**:437-443.
- 662 36. **González-Escalona N, Martínez-Urtaza J, Romero J, Espejo RT, Jaykus L-A,**  
663 **DePaola A.** 2008. Determination of molecular phylogenetics of *Vibrio parahaemolyticus*  
664 strains by multilocus sequence typing. *J Bacteriol* **190**:2831-2840.
- 665 37. **Ellingsen BA, Olsen JS, Granum PE, Rorvik LM, González-Escalona N.** 2013.  
666 Genetic characterization of *trh* positive *Vibrio* spp. isolated from Norway. *Front Cell*  
667 *Infect Microbiol* **3**:107.
- 668 38. **Chen Y, Dai J, Morris JG, Johnson JA.** 2010. Genetic analysis of the capsule  
669 polysaccharide (K antigen) and exopolysaccharide genes in pandemic *Vibrio*  
670 *parahaemolyticus* O3: K6. *BMC Microbiol* **10**:1.
- 671 39. **Meibom KL, Blokesch M, Dolganov NA, Wu C-Y, Schoolnik GK.** 2005. Chitin  
672 induces natural competence in *Vibrio cholerae*. *Science* **310**:1824-1827.
- 673 40. **Takemura AF, Chien DM, Polz MF.** 2014. Associations and dynamics of *Vibrionaceae*  
674 in the environment, from the genus to the population level. *Front Microbiol* **5**:38.
- 675 41. **Lovell CR.** 2017. Ecological fitness and virulence features of *Vibrio parahaemolyticus* in  
676 estuarine environments. *Appl Microbiol Biotechnol* **101**:1781-1794.
- 677 42. **Johnson CN.** 2013. Fitness factors in vibrios: a mini-review. *Microb Ecol* **65**:826-851.
- 678 43. **Matz C, Nouri B, McCarter L, Martínez-Urtaza J.** 2011. Acquired type III secretion  
679 system determines environmental fitness of epidemic *Vibrio parahaemolyticus* in the  
680 interaction with bacterivorous protists. *PloS one* **6**:e20275.

- 681 44. **Nishibuchi M, Kaper JB.** 1985. Nucleotide sequence of the thermostable direct  
682 hemolysin gene of *Vibrio parahaemolyticus*. J Bacteriol **162**:558-564.
- 683 45. **Baym M, Kryazhimskiy S, Lieberman TD, Chung H, Desai MM, Kishony R.** 2015.  
684 Inexpensive multiplexed library preparation for megabase-sized genomes. PloS one  
685 **10**:e0128036.
- 686 46. **Tritt A, Eisen JA, Facciotti MT, Darling AE.** 2012. A5. An integrated pipeline for *de*  
687 *novo* assembly of microbial genomes. PLoS ONE **7**:e42304.
- 688 47. **Seemann T.** 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics.  
689 **30**:2068-9
- 690 48. **Inouye M, Conway TC, Zobel J, Holt KE.** 2012. Short read sequence typing (SRST):  
691 multi-locus sequence types from short reads. BMC Genomics **13**:338.
- 692 49. **Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA,**  
693 **Zeng Q, Wortman J, Young SK.** 2014. Pilon: an integrated tool for comprehensive  
694 microbial variant detection and genome assembly improvement. PloS one **9**:e112963.
- 695 50. **Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden**  
696 **TL.** 2009. BLAST+: architecture and applications. BMC Bioinformatics **10**:421.
- 697 51. **Sullivan MJ, Petty NK, Beatson SA.** 2011. Easyfig: a genome comparison visualizer.  
698 Bioinformatics **27**:1009-1010.
- 699 52. **Gardner SN, Slezak T, Hall BG.** 2015. kSNP3. 0: SNP detection and phylogenetic  
700 analysis of genomes without genome alignment or reference genome. Bioinformatics  
701 **31**:2877-8.
- 702 53. **Stamatakis A.** 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic  
703 analyses with thousands of taxa and mixed models. Bioinformatics **22**:2688-2690.

- 704 54. **Darling AC, Mau B, Blattner FR, Perna NT.** 2004. Mauve: multiple alignment of  
705 conserved genomic sequence with rearrangements. *Genome Res* **14**:1394-1403.
- 706 55. **Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M,**  
707 **Falush D, Keane JA, Parkhill J.** 2015. Roary: rapid large-scale prokaryote pan genome  
708 analysis. *Bioinformatics* **31**:3691-3693.
- 709 56. **Sahl JW, Lemmer D, Travis J, Schupp J, Gillece J, Aziz M, Driebe E, Drees K,**  
710 **Hicks N, Williamson C.** 2016. The Northern Arizona SNP Pipeline (NASP): accurate,  
711 flexible, and rapid identification of SNPs in WGS datasets. *Microb Genom.* **2**:e000074
- 712 57. **Whistler CA, Hall JA, Xu F, Ilyas S, Siwakoti P, Cooper VS, Jones SH.** 2015. Use of  
713 Whole-Genome Phylogeny and Comparisons for Development of a Multiplex PCR Assay  
714 To Identify Sequence Type 36 *Vibrio parahaemolyticus*. *J Clin Microbiol* **53**:1864-1872.
- 715 58. **Jolley KA, Chan M-S, Maiden MC.** 2004. mlstdbNet—distributed multi-locus sequence  
716 typing (MLST) databases. *BMC Bioinformatics* **5**:86.
- 717 59. **Alikhan N-F, Petty NK, Zakour NLB, Beatson SA.** 2011. BLAST Ring Image  
718 Generator (BRIG): simple prokaryote genome comparisons. *BMC Genomics* **12**:402  
719  
720

721 Table 1: Clinical and environmental prevalence of emergent Northeast US *V. parahaemolyticus*  
 722 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>e</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

723 <sup>a</sup> Some clinical isolates had insufficient sequencing coverage to determine sequence type and included  
 724 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  
 725 which were from wound infections. Two wound infection isolates lacking hemolysins were of known  
 726 sequence types and are not listed above.

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

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

730 <sup>d</sup> Presence of the VPαIγ architecture was determined by PacBio genome sequencing of isolate MAVP-Q  
 731 and MAVP-26, whereas for other isolates, identification of VPαI type was determined through illumina  
 732 genome sequencing, PCR amplification and Sanger sequencing.

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

735  
 736

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

743

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

750

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

755

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

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

762

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



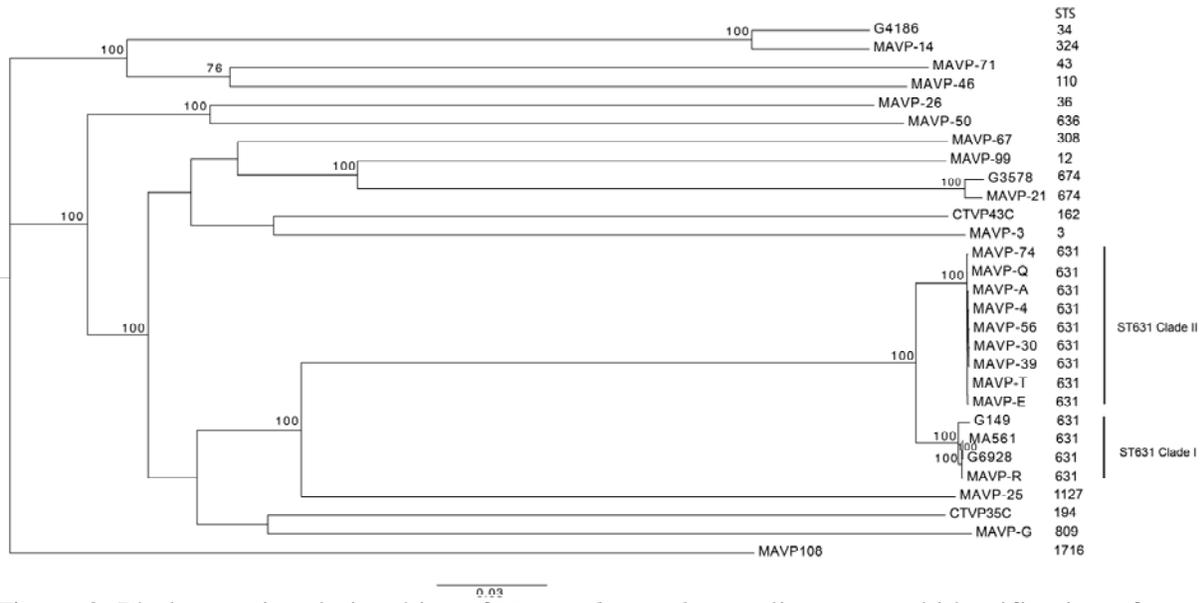


Figure 2. Phylogenetic relationships of *V. parahaemolyticus* lineages and identification of distinct ST631 clades. An ML phylogeny of representative *V. parahaemolyticus* genomes of clinical strains causing two or more infections was built on whole genome SNPs identified by reference-free comparisons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



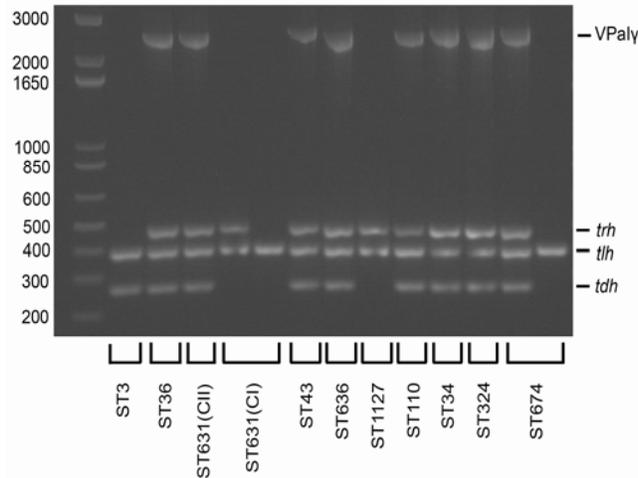


Figure 5. Distribution of VPaly in emergent pathogen lineages. The presence of *tdh*, *trh* and VPaly along with positive control *tlh* was determined by PCR amplification using gene-specific primers and visualized on a 1.2% agarose gel. The order from left to right is 1kb+ ladder, ST3 (MAVP-C), ST36 (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.