

1 **Evaluation of whole genome sequencing for the identification and typing of *Vibrio***

2 ***cholerae***

3

4

5 David R. Greig¹, Ulf Schafer¹, Sophie Octavia^{2,3}, Ebony Hunter^{1,4}, Marie A. Chattaway¹,

6 Timothy J. Dallman¹, Claire Jenkins*¹

7

8 ¹National Infection Services, Public Health England, 61 Colindale Avenue, London NW9 5HT,

9 UK.

10 ²National Public Health Laboratory, Ministry of Health, Singapore

11 ³School of Biotechnology and Biomolecular Sciences, University of New South Wales, Sydney,

12 Australia

13 ⁴School of Pharmacy & Biomolecular Sciences, University of Brighton, Sussex, UK

14

15 * Corresponding author's details:

16 Address: Gastrointestinal Bacteria Reference Unit,

17 Public Health England,

18 61 Colindale Avenue,

19 London, UK

20 NW9 5HT

21

22 Running title – WGS for *Vibrio cholerae*

23

24 **Abstract**

25 Epidemiological and microbiological data on *Vibrio cholerae* isolated between 2004 and
26 2017 (n=836) and held in the Public Health England culture archive were reviewed. The
27 traditional biochemical species identification and serological typing results were compared
28 with the genome derived species identification and serotype for a sub-set of isolates
29 (n=152). Of the 836 isolates, 750 (89.7%) were from faecal specimens, 206 (24.6%)
30 belonged to serogroup O1 and seven (0.8%) were serogroup O139, and 792 (94.7%) isolates
31 from patients reporting recent travel abroad, most commonly to India (n=209) and Pakistan
32 (n=104). Of the 152 isolates of *V. cholerae* speciated by kmer identification, 149 (98.1%)
33 were concordant with the traditional biochemical approach. Traditional serotyping results
34 were 100% concordant with the whole genome sequencing (WGS) analysis for identification
35 of serogroups O1 and O139 and Classical and El Tor biotypes. *ctxA* was detected in all
36 isolates of *V. cholerae* O1 El Tor and O139 belonging to sequence type (ST) 69, and in *V.*
37 *cholerae* O1 Classical variants belonging to ST73. A phylogeny of isolates belonging to ST69
38 from UK travellers clustered geographically, with isolates from India and Pakistan located on
39 separate branches. Moving forward, WGS data from UK travellers will contribute to global
40 surveillance programs, and the monitoring of emerging threats to public health and the
41 global dissemination of pathogenic lineages. At the national level, these WGS data will
42 inform the timely reinforcement of direct public health messaging to travellers and mitigate
43 the impact of imported infections and the associated risks to public health.

44 **Introduction**

45 Cholera is an acute diarrhoeal disease that can kill within hours if left untreated. Patients
46 present with the passing of voluminous rice water stools leading to severe dehydration (1).
47 If hydration and electrolyte therapy is not quickly initiated, symptoms can rapidly progress
48 to hypovolemic shock, acidosis and death. Inadequate access to clean water and sanitation
49 facilities is a driver of transmission, and outbreaks are common among displaced
50 populations living in overcrowded conditions (2).

51

52 The bacterial pathogen responsible for the disease is *Vibrio cholerae*. *V. cholerae*
53 serogroups O1 and O139 are regarded as pandemic strains and harbour the *ctx* genes
54 associated with the production of cholera toxin (3). *ctx* has also been detected in a limited
55 number of other serogroups (4). Serogroup O1 can be divided into two biotypes, Classical
56 and El Tor. There are over 200 different lipopolysaccharide 'O' antigens or serogroups of *V.*
57 *cholerae*. The non-O1, non-O139 serogroups are associated with a milder form of
58 gastroenteritis, septicaemia and other extra-intestinal infections (1, 3).

59

60 Seven cholera pandemics have occurred throughout the 19th and 20th centuries. The
61 seventh (and current) pandemic began in the Bay of Bengal and has spread to Africa and
62 South America in at least three independent but overlapping waves of transmission (5). The
63 fifth and sixth pandemics were caused by the *V. cholerae* serogroup O1 Classical biotype
64 while the seventh pandemic was caused by serogroup O1 biotype El Tor. In 1992, *V.*
65 *cholerae* serogroup O139 caused a large epidemic in Bangladesh and India (6), however *V.*
66 *cholerae* O1 El Tor persists as the most commonly isolated *ctx*-positive serotype/biotype. *V.*
67 *cholerae* is endemic across Africa, Latin America and Asia resulting in a large healthcare
68 burden in developing countries (7-9). The World Health Organisation states that there are
69 1.3 - 4 Million estimated cases and 21,000- 147,000 estimated deaths annually (7).

70

71 The UK Standards for Microbiology Investigations Investigation of Faecal Specimens for
72 Enteric Pathogens recommends testing of faecal specimens for *V. cholerae* in cases of
73 suspected cholera, seafood consumption, and/or recent travel (within 2-3 weeks) to
74 countries where cholera is endemic ([https://www.gov.uk/government/publications/smi-b-
75 30-investigation-of-faecal-specimens-for-enteric-pathogens](https://www.gov.uk/government/publications/smi-b-30-investigation-of-faecal-specimens-for-enteric-pathogens)). Consequently, the true
76 incidence of domestically acquired *V. cholerae* in the UK is unknown, and almost all isolates
77 of enteric origin are associated with travellers' diarrhoea.

78

79 In 2015, Public Health England (PHE) implemented whole genome sequencing (WGS) for the
80 routine surveillance of the more common gastrointestinal pathogens including *E.coli*,
81 *Salmonella*, *Campylobacter*, *Shigella* and *Listeria* species (10-12). The aim of the study was
82 to review the historical PHE data on isolates of *V. cholerae* held in the PHE culture archives,
83 compare the results of the traditional biochemical and serological methods with the analysis
84 of WGS data for a sub-set of isolates, and assess the impact of implementing WGS for the
85 public health surveillance of *V. cholerae*.

86

87 **Methods**

88 *Epidemiological data*

89 All isolates of *V. cholerae* from human cases resident in England submitted to the
90 Gastrointestinal Bacteria Reference Unit (GBRU) by local hospital laboratories between 2004
91 and 2017 were reviewed. Patient information including, sex, age and recent travel, was
92 collected from laboratory request forms upon submission and stored in the Gastro Data
93 Warehouse (GDW), an in-house PHE database for storing and linking patient demographic
94 and microbiological typing data. Data on symptoms were limited stating only that the

95 patient had either gastrointestinal symptoms or an extra-intestinal infection. There was no
96 data on severity of symptoms or patient outcome.

97

98 *Bacterial culture and traditional biochemistry and serology*

99 Cultures were stored on cryobeads at -40°C or in nutrient agar stabs. For each sample, one
100 cryobead was taken and inoculated into 20ml 3% NaCl peptone water and incubated at 37°C
101 for 18 hours, shaking at 80rpm. Cultures were plated out from either nutrient agar slopes or
102 3% NaCl peptone water onto Blood agar, MacConkey agar with salt (NaCl 1%), Thiosulphate-
103 citrate-bile salts (TCBS) agar and cystine lactose electrolyte deficient (CLED) agar and
104 incubated at 37°C for 18 hours.

105

106 Biochemical identification was performed following inoculation onto a panel of substrates.

107 Utilisation of the substrate was identified by a colour changes or gas production within the

108 media. All positive and negative reactions were compared to a known reference panel of

109 results to give a final identification. Isolates of *V. cholerae* were agglutinated with antisera

110 raised to O1 (Ogawa and Inaba) and O139 (Bengal) antisera to determine the serogroup.

111 The Classical and El-Tor biotypes were differentiated by the Voges-Proskauer (VP) test

112 (Classical, negative; El-Tor, positive) and haemolysis on blood agar (Classical, non-

113 haemolytic; El-Tor, haemolytic).

114

115 *Whole genome sequencing analysis*

116 All viable cultures of *V. cholerae* submitted to GBRU between January 2015 and March 2018

117 were sequenced (n=152). Genomic DNA was extracted, fragmented and tagged for

118 multiplexing with Nextera XT DNA Sample Preparation Kits (Illumina) and sequenced using

119 the Illumina HiSeq 2500 at PHE. FASTQ reads were quality trimmed using Trimmomatic

120 (v0.36) (13) with bases removed from the trailing end that fell below a PHRED score of 30. If

121 the read length post trimming was less than 50, the read and its pair were discarded using
122 Trimmomatic. FASTQ reads from all sequences in this study can be found at the PHE
123 Pathogens BioProject at the National Center for Biotechnology Information (PRJNA438219).
124
125 A kmer (a short string of DNA of length k; in this method k=18) based approach was used to
126 confirm the identity of the sample before organism specific algorithms were applied
127 (<https://github.com/phe-bioinformatics/kmerid>) (14). Reference genomes (n=1781) in 59
128 bacterial genera comprising the majority of human pathogens, commensal bacteria and
129 common contaminants were downloaded from
130 <ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq/bacteria>. The kmer algorithm compared each
131 sample to representative genomes in these 59 bacterial genera and returned the most
132 similar genome together with a similarity estimate.
133
134 Sequence Type (ST) assignment was performed using a modified version of SRST using the
135 MLST database described by Tewolde *et al.* 2016 (15). The MOST software (for MLST) is
136 available at <https://github.com/phe-bioinformatics/MOST>. Any MLST gene sequences that
137 did not match the existing alleles were submitted to pubMLST
138 (<https://pubmlst.org/vcholerae/>) for a new allelic type assignment. Similarly, new allelic
139 profiles were also submitted to the database for a new sequence type (ST) assignment.
140
141 For the isolates belonging to clonal complex (CC) 69, high quality Illumina reads were
142 mapped to a SPADES v3.5.0 *de novo* assembly of the *V. cholerae* reference genomes NC-
143 002505.1 and NC-002506.1, using BWA-MEM v0.7.3 and Samtools v1.1 (16, 17). Single
144 Nucleotide Polymorphisms (SNPs) were identified using GATK v2.6.5 (18) in unified
145 genotyper mode. Core genome positions, defined as those present in the reference genome
146 and at least 80% of the isolates, that had a high quality SNP (>90% consensus, minimum

147 depth 10x, MQ >= 30) in at least one isolate were extracted using SnapperDB v0.2.5 and
148 processed though Gubbins v2.0.0 to account and suppress recombination within the input to
149 RAxML v8.1.17 (19).

150

151 Using the *GeneFinder* tool (Doumith, unpublished), FASTQ reads were mapped to the
152 virulence regulator gene, *toxR* (Genbank accession: KF498634.1), the cholera toxin gene *ctxA*
153 (Genbank accession: AF463401.1), *wbeO1*, and *wbfO139* (Genbank accessions: KC152957.1
154 and AB012956.1) encoding the somatic O antigens O1 and O139, *tcpA* classical and *tcpA* El
155 Tor gene sequences (Genbank accessions: M33514.1 and KP187623.1) using Bowtie 2 (20).
156 The best match to each target was reported with metrics including coverage, depth and
157 nucleotide similarity in XML format for quality assessment. *toxR* is found in all isolates of *V.*
158 *cholerae* and are regarded as a marker for species identification, *ctxA* encoding cholera toxin
159 is associated with *V. cholerae* O1 and O139 and is a marker for the pandemic lineages (21).
160 Variants of *tcpA* can be used to identify the Classical and El Tor biotypes (21). For *in silico*
161 predictions, only results that matched to a gene determinant at >80% nucleotide identity
162 over >80% target gene length were accepted.

163

164 **Results**

165 *Review of the historical data*

166 Between January 2014 and December 2017, 836 isolates of *V. cholerae* from human cases
167 resident in England were submitted to GBRU by local hospital laboratories. On average, the
168 number of isolates per year was 60, with the lowest number of isolates being reported in
169 2013 (n=29) and the highest number was reported in 2007 (n=80) (Figure 1). Of the 836
170 isolates, 206 (24.6%) belonged to serogroup O1 and seven were serogroup O139 (0.8%), and
171 750 were from faecal specimens, six were from blood cultures, two were from ear swabs
172 and two were from eye swabs. No clinical data was available for the remaining 76 isolates.

173

174 Gender and age data was available for 828/836 and 773/836 cases, respectively. There were
175 424/836 males (50.7%) and 404/836 females (48.3%), and 685/836 (81.9%) were adults
176 (aged 16 years or older) and 88 (10.5%) were children (<16 years old). Travel data was
177 available for 796/836 cases, of which 792 reported recent travel abroad (less than 7 days of
178 onset of symptoms). For the cases infected with *V. cholerae* non-O1, non-O139, the most
179 common travel destinations were India (n=140), Kenya (n=57), Thailand (n=40) and Egypt
180 (n=36). The most commonly reported destinations of cases of *V. cholerae* O1 were Pakistan
181 (n=72) and India (n=69). The six cases of *V. cholerae* O139 had travelled to Thailand (n=2),
182 China, India, Jordan and Pakistan. The four cases who stated they had not recently travelled
183 abroad, all had *V. cholerae* non-O1, non-O139 isolates from extra-intestinal sites (blood
184 cultures, n=2; eye swab, n=1; ear swab, n=1).

185

186 *Whole genome sequencing*

187 One hundred and fifty-two isolates of *V. cholerae* were sequenced including those belonging
188 to serogroup O1 (n=47), serogroup O139 (n=7), those designated serogroup non-O1, non-
189 O139 (n=98) (Supplementary Table 1). One hundred and thirty-seven were isolates from
190 human cases, of which 132 were from faecal specimens from hospitalised or community
191 cases with symptoms of gastrointestinal disease, three isolates were from ear swabs, one
192 was from an eye swab and one was from a blood culture from a patient with acute
193 cholecystitis. The remaining 15 isolates were from animals (n=4), food (n=1), environmental
194 samples (n=3) or were isolates from the National Collection of Type Cultures (n=7)
195 (Supplementary Table 1).

196

197 *Kmer identification*

198 Of the 152 isolates of *V. cholerae* speciated using the kmer identification approach, 149
199 (98.1%) were concordant with the traditional biochemical identification. The kmer method
200 failed to identify three unusual external quality assessment isolates from an obscure
201 environmental source, previously identified as *V. cholerae*. These isolates were *V. cholerae*,
202 however, the similarity of the sequences to the *V. cholerae* reference sequences in the kmer
203 ID database was below the acceptable threshold (80% similarity) required to confirm the
204 identification.

205

206 *Use of Genefinder for serotyping and detection of virulence genes*

207 The *toxR* gene was detected in 144/152 isolates of *V. cholerae* (Table 1). Eight isolates
208 identified as *V. cholerae* by traditional biochemical tests, were negative for *toxR*. Further
209 analysis of the sequences data showed the sequence coverage/similarity of *toxR* for the
210 discrepant isolates fell just below the 80% threshold (74% and 77%) (Table 1). However, all
211 eight isolates were identified as *V. cholerae* by the kmer approach.

212

213 Traditional biochemistry and serotyping results were 100% concordant with the WGS
214 analysis for identification of O1 and O139 and Classical and El Tor biotypes. Of the 37
215 isolates of *V. cholerae* O1 El Tor, 33 had *wbeO1*, *tcpA* El Tor variant and *ctxA* and four had
216 had *wbeO1*, *tcpA* El Tor variant without *ctxA* (Table 1). There were five isolates that
217 belonged to *V. cholerae* serogroup O1 that were negative for *tcpA* El Tor variant and *ctxA*,
218 and one that was negative for *tcpA* El Tor variant but had *ctxA* (Table 1). There were only
219 three Classical variant strains in the study, as the current pandemic is caused by *V. cholerae*
220 O1 El Tor (7-9, 22, 23). All three *V. cholerae* O1 classical strains had *wbeO1* and the *tcpA*
221 classical variant gene. *V. cholerae* O139 was also rare in this dataset (24). Although all
222 seven isolates of *V. cholerae* O139 had *wbfO139*, only the NCTC strain had the *tcpA* El Tor
223 variant and *ctxA* (Table 1).

224

225 *ctxA* was detected in all isolates of *V. cholerae* O1 El Tor and O139 belonging to ST69 (25)
226 and *V. cholerae* O1 Classical variants belonging to ST73 (22). Four isolates of *V. cholerae* O1
227 were negative for *ctxA*, and the six recently isolated of *V. cholerae* O139, were *ctxA*-
228 negative.

229

230 *Sequence typing*

231 Sequence typing data was available for 152 isolates. The *V. cholerae* O139 El Tor -positive
232 isolate and the 34 isolates of *V. cholerae* O1 El Tor *ctxA*-positive isolates belonged to ST69.
233 The four *V. cholerae* O1 El Tor *ctxA*-negative isolates belonged to ST75, ST169 and ST579
234 (n=2) and all fell within CC69. Previously studies have suggested that the emergence and
235 potential spread of ST75 may pose significant threat to public health (26). Epidemiological
236 surveillance is required to further understand the epidemic potential of *ctxA*-negative STs
237 that are part of CC69. The *V. cholerae* O1 classical isolate belonged to ST73.

238

239 The six isolates of *V. cholerae* O1 without the *tcpA* El Tor variant were ST167, ST521 (n=2)
240 and ST551 (n=2) and ST611. There were six isolates that had the O139 antigen but were
241 negative for the *tcpA* El Tor variant gene and *ctxA*, and these belonged to ST163, ST527,
242 ST529, ST544, ST568 and ST586. All *V. cholerae* O1 isolates, regardless of the presence or
243 absence of *tcpA* El Tor variant or *ctxA*, belonged to the CC69 cluster and those without the
244 *tcpA* El Tor variant gene were dispersed across the population (Figure 1).

245

246 The remaining 95 isolates of *V. cholerae* non-O1, non-O139 (n=95) and *V. cholerae* O139
247 (n=3), belonged to over 70 different STs (Supplementary Table). There was only one major
248 cluster among the *V. cholerae* non-O1, non-O139 isolates, designated CC558. Isolates
249 belonging to this cluster were geographically dispersed (Figure 1).

250

251 *SNP typing*

252 As previously described, the pandemic *V. cholerae* O1 and O139 El Tor *ctxA* strains all
253 belonged to ST69, whereas the Classical biotype and *ctxA*-negative strains of *V. cholerae* O1
254 belonged to other STs within CC69. A phylogeny of ST69, the pandemic lineage, was
255 constructed comprising isolates from this study and sequences available in public databases
256 (Supplementary Figure). The isolates from UK travellers clustered geographically with those
257 returning from India located on the same branch, and those reporting recent travel to
258 Pakistan clustered on a separate branch. Further analysis based on single nucleotide
259 polymorphisms in the core genome compared to a reference strain may be performed for
260 outbreak detection and source attribution where the incidence of the current *V. cholerae* O1
261 El Tor pandemic lineage (ST69) is high (27, 28)

262

263 **Discussion**

264 Historically, traditional biochemistry, biotyping, phage typing and serology results were
265 useful for confirming identification at the species level, typing of serogroups O1 and O139
266 and for identifying the Classical and El Tor variants. Isolates belonging to serogroups O1 and
267 O139 were assumed to belong to the pandemic lineages and have the potential to cause
268 cholera. In this study, the review of the historical GBRU data revealed that just under a
269 quarter of the isolates of *V. cholerae* belonged to serogroup O1 and *V. cholerae* O139 was
270 rarely detected (22). Due to limited resources, neither serotyping of the non-O1, non-O139
271 serogroups, nor molecular typing of any serogroup, were performed at GBRU. Therefore,
272 prior to the implementation of WGS, it was not possible to monitor trends in emerging
273 pathogenic lineages or gain insight into modes of transmission for this important
274 gastrointestinal pathogen.

275

276

277 Previous studies have shown that MLST data is an accurate, robust, reliable, high throughput
278 typing method that is well suited to routine public health surveillance (11, 12, 15). For *V.*
279 *cholerae*, MLST provides insight on the true evolutionary relationship between isolates, as
280 well as a framework for fine level typing for public health surveillance (29-32). Using the STs
281 derived from the genome data, we were able to analyse the population structure of all
282 isolates of *V. cholerae* submitted to GBRU for the first time. As previous studies have
283 shown, the population structure of the non-O1 and non-O139 serotypes was diverse (31).
284 Currently, the correlation of ST with geography in our dataset is hindered by the limited size
285 of the dataset. However, moving forward this unprecedented level of strain discrimination
286 available for all isolates of *V. cholerae* submitted to GBRU will enhance our understanding of
287 global transmission and emerging threats to public health, for the pandemic strains
288 belonging to CC69, and the non-O1 serogroup lineages.

289

290 A review of the data on cases of travellers' diarrhoea caused by *V. cholerae* held by GBRU
291 showed that travel histories, including the country visited, were complete for 95.2% of
292 cases. Therefore, these data have the potential to be a useful public health resource for
293 global surveillance, enabling us to track the emergence and dissemination of specific
294 lineages on a global scale (5, 8, 33). Furthermore, at the national level, sharing of WGS data
295 linked to these cases could result in the timely reinforcement of direct public health
296 messaging to travellers, in order to reduce the number of imported infections and mitigate
297 the impact of imported infections and associated risks to public health (34).

298

299 The consequence of humanitarian crises, such as the disruption of water and sanitation
300 systems and the displacement of populations to overcrowded camps, increases the risk of
301 the transmission and outbreaks of cholera (2). However, robust global monitoring of *V.*

302 *cholerae* is hindered by the limitations of the surveillance systems in countries where people
303 are most at risk. The World Health Organisation recommends that cholera surveillance
304 should be part of an integrated disease surveillance system that includes feedback at the
305 local level and information sharing at the global level
306 (<http://www.who.int/mediacentre/factsheets/fs107/en/>).

307

308 Traditional biochemistry and serotyping results were concordant with the WGS analysis for
309 identification of *V. cholerae* O1, serotyping and biotyping of O1 and O139 serogroups.
310 Moreover, using the WGS approach species level identification, serotyping, biotyping,
311 presence of cholera toxin, ST and SNP typing of CC69, can all be derived from a single
312 process work flow. WGS data may also be interrogated for additional virulence factors, and
313 antimicrobial resistance determinants. The genomic data of all *V. cholerae* sequenced at
314 PHE are publically released into the NCBI BioProject PRJNA438219 in order to facilitate
315 public health surveillance, and monitoring of the global transmission of the pandemic
316 lineages, by the international scientific community.

317

318

319 **Funding**

320 This study was funded by Public Health England and supported by the National Institute for
321 Health Research Health Protection Research Unit in Gastrointestinal Infections (#109524).

322 The views expressed are those of the author(s) and not necessarily those of the NHS, the
323 NIHR, the Department of Health or Public Health England.

324

325

326

327 **References**

- 328 1. **Harris JB, LaRocque RC, Qadri F, Ryan ET, Calderwood SB.** 2012. Cholera. *Lancet*
329 **379**:2466-76.
- 330 2. **Jutla A, Khan R, Colwell R.** 2017. Natural Disasters and Cholera Outbreaks: Current
331 Understanding and Future Outlook. *Curr Environ Health Rep* **4**:99-107.
- 332 3. **Islam MT, Alam M, Boucher Y.** 2017. Emergence, ecology and dispersal of the
333 pandemic generating *Vibrio cholerae* lineage. *Int Microbiol* **20**:106-115.
- 334 4. **Crowe SJ, Newton AE, Gould LH, Parsons MB, Stroika S, Bopp CA, Freeman M,**
335 **Greene K, Mahon BE.** 2016. Vibriosis, not cholera: toxigenic *Vibrio cholerae* non-
336 O1, non-O139 infections in the United States, 1984-2014. *Epidemiol Infect*
337 **144**:3335-3341.
- 338 5. **Mutreja A, Kim DW, Thomson NR, Connor TR, Lee JH, Kariuki S, Croucher NJ, Choi**
339 **SY, Harris SR, Lebens M, Niyogi SK, Kim EJ, Ramamurthy T, Chun J, Wood JL,**
340 **Clemens JD, Czerkinsky C, Nair GB, Holmgren J, Parkhill J, Dougan G.** 2011.
341 Evidence for several waves of global transmission in the seventh cholera pandemic.
342 *Nature* **477**:462-5.
- 343 6. **Albert MJ.** 1996. Epidemiology & molecular biology of *Vibrio cholerae* O139 Bengal.
344 *Indian J Med Res* **104**:14-27.
- 345 7. **Clemens JD, Nair GB, Ahmed T, Qadri F, Holmgren J.** 2017. Cholera. *Lancet*
346 **390**:1539-1549.
- 347 8. **Domman D, Quilici ML, Dorman MJ, Njamkepo E, Mutreja A, Mather AE, Delgado**
348 **G, Morales-Espinosa R, Grimont PAD, Lizárraga-Partida ML, Bouchier C, Aanensen**
349 **DM, Kuri-Morales P, Tarr CL, Dougan G, Parkhill J, Campos J, Cravioto A, Weill**
350 **FX, Thomson NR.** 2017. Integrated view of *Vibrio cholerae* in the Americas. *Science*
351 **358**:789-793.

- 352 9. **Weill FX, Domman D, Njamkepo E, Tarr C, Rauzier J, Fawal N, Keddy KH, Salje H,**
353 **Moore S, Mukhopadhyay AK, Bercion R, Luquero FJ, Ngandjio A, Dosso M,**
354 **Monakhova E, Garin B, Bouchier C, Pazzani C, Mutreja A, Grunow R, Sidikou F,**
355 **Bonte L, Breurec S, Damian M, Njanpop-Lafourcade BM, Sapriel G, Page AL, Hamze**
356 **M, Henkens M, Chowdhury G, Mengel M, Koeck JL, Fournier JM, Dougan G,**
357 **Grimont PAD, Parkhill J, Holt KE, Piarroux R, Ramamurthy T, Quilici ML, Thomson**
358 **NR.** 2017. Genomic history of the seventh pandemic of cholera in Africa. *Science*
359 **358**:785-789.
- 360 10. **Dallman TJ, Byrne L, Ashton PM, Cowley LA, Perry NT, Adak G, Petrovska L, Ellis**
361 **RJ, Elson R, Underwood A, Green J, Hanage WP, Jenkins C, Grant K, Wain J.** 2015.
362 Whole-genome sequencing for national surveillance of Shiga toxin-producing
363 *Escherichia coli* O157. *Clin Infect Dis* **61**:305-12.
- 364 11. **Ashton PM, Nair S, Peters TM, Bale JA, Powell DG, Painset A, Tewolde R, Schaefer**
365 **U, Jenkins C, Dallman TJ, de Pinna EM, Grant KA; Salmonella Whole Genome**
366 **Sequencing Implementation Group.** 2016. Identification of Salmonella for public
367 health surveillance using whole genome sequencing. *PeerJ* **4**:e1752.
- 368 12. **Chattaway MA, Greig DR, Gentle A, Hartman HB, Dallman TJ, Jenkins C.** 2017.
369 Whole-Genome Sequencing for National Surveillance of *Shigella flexneri*. *Front*
370 *Microbiol* **8**:1700.
- 371 13. **Bolger AM, Lohse M, Usadel B.** 2014. Trimmomatic: A flexible trimmer for Illumina
372 Sequence Data. *Bioinformatics* **30**:2114-20.
- 373 14. **Chattaway MA, Schaefer U, Tewolde R, Dallman TJ, Jenkins C.** 2017. Identification
374 of *Escherichia coli* and *Shigella* species from Whole-Genome Sequences. *J Clin*
375 *Microbiol* **55**:616-623.
- 376
377

- 378 15. **Tewolde R, Dallman T, Schaefer U, Sheppard CL, Ashton P, Pichon B, Ellington M,**
379 **Swift C, Green J, Underwood A.** 2016. MOST: a modified MLST typing tool based on
380 short read sequencing. *PeerJ* **4**:e2308.
- 381 16. **Li H, Durbin R.** 2010. Fast and accurate long-read alignment with Burrows-Wheeler
382 transform. *Bioinformatics* **26**:589–595.
- 383 17. **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,**
384 **Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G,**
385 **Alekseyev MA, Pevzner PA.** 2012. SPAdes: a new genome assembly algorithm and
386 its applications to single-cell sequencing. *J Comput Biol* **19**:455–477
- 387 18. **McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella**
388 **K, Altshuler D, Gabriel S, Daly M, DePristo MA.** 2010 The Genome Analysis Toolkit:
389 a MapReduce framework for analyzing next-generation DNA sequencing data.
390 *Genome Research* **20**:1297–1303.
- 391 19. **Stamatakis A.** 2014 RAxML version 8: a tool for phylogenetic analysis and post-
392 analysis of large phylogenies. *Bioinformatics* **30**:1312–1313.
- 393 20. **Langmead B, Salzberg SL.** 2012. Fast gapped-read alignment with Bowtie 2. *Nat*
394 *Methods* **9**:357-359.
- 395 21. **Greig DR, Hickey TJ, Boxall MD, Begum H, Gentle A, Jenkins C, Chattaway MA.**
396 2018. A real-time multiplex PCR for the identification and typing of *Vibrio cholerae*.
397 *Diagn Microbiol Infect Dis* **90**:171-176.
- 398 22. **Mukhopadhyay AK, Takeda Y, Balakrish Nair G.** 2014. Cholera outbreaks in the El
399 Tor biotype era and the impact of the new El Tor variants. *Curr Top Microbiol*
400 *Immunol* **379**:17-47.
- 401 23. **Siddique AK, Cash R.** 2014. Cholera outbreaks in the classical biotype era. *Curr Top*
402 *Microbiol Immunol* **379**:1-16.

- 403 24. **Ghosh R, Sharma NC, Halder K, Bhadra RK, Chowdhury G, Pazhani GP, Shinoda S,**
404 **Mukhopadhyay AK, Nair GB, Ramamurthy T.** 2016. Phenotypic and Genetic
405 Heterogeneity in *Vibrio cholerae* O139 Isolated from Cholera Cases in Delhi, India
406 during 2001-2006. *Front Microbiol* **7**:1250.
- 407 25. **Anandan S, Devanga Ragupathi NK, Muthuirulandi Sethuvel DP, Thangamani S,**
408 **Veeraraghavan B.** 2017. Prevailing clone (ST69) of *Vibrio cholerae* O139 in India over
409 10 years. *Gut Pathog.* **9**:60.
- 410 26. **Luo Y, Octavia S, Jin D, Ye J, Miao Z, Jiang T, Xia S, Lan R.** 2016. US Gulf-like
411 toxigenic O1 *Vibrio cholerae* causing sporadic cholera outbreaks in China. *J Infect*
412 **72**:564-72.
- 413 27. **Ramamurthy T, Sharma NC.** 2014. Cholera outbreaks in India. *Curr Top Microbiol*
414 *Immunol* **379**:49-85.
- 415 28. **Shah MA, Mutreja A, Thomson N, Baker S, Parkhill J, Dougan G, Bokhari H, Wren**
416 **BW.** 2014. Genomic epidemiology of *Vibrio cholerae* O1 associated with
417 floods, Pakistan, 2010. *Emerg Infect Dis* **20**:13-20.
- 418 29. **Kotetishvili M, Stine OC, Chen Y, Kreger A, Sulakvelidze A, Sozhamannan S, Morris**
419 **JG Jr. 2003.** Multilocus sequence typing has better discriminatory ability for
420 typing *Vibrio cholerae* than does pulsed-field gel electrophoresis and provides a
421 measure of phylogenetic relatedness. *J Clin Microbiol* **41**:2191-6
- 422 30. **Lam C, Octavia S, Reeves PR, Lan R.** 2012. Multi-locus variable number tandem
423 repeat analysis of 7th pandemic *Vibrio cholerae*. *BMC Microbiol* **12**:82.
- 424 31. **Octavia S, Salim A, Kurniawan J, Lam C, Leung Q, Ahsan S, Reeves PR, Nair GB, Lan**
425 **R.** 2013. Population structure and evolution of non-O1/non-O139 *Vibrio cholerae* by
426 multilocus sequence typing. *PLoS One* **8**:e65342.
- 427 32. **Siriphap A, Leekitcharoenphon P, Kaas RS, Theethakaew C, Aarestrup FM,**
428 **Sutheinkul O, Hendriksen RS.** 2017. Characterization and genetic variation of *Vibrio*

- 429 *cholerae* isolated from clinical and environmental sources in Thailand. PLoS One
430 **12**:e0169324.
- 431 33. **Chowdhury FR, Nur Z, Hassan N, von Seidlein L, Dunachie S.** 2017. Pandemics,
432 pathogenicity and changing molecular epidemiology of cholera in the era of global
433 warming. *Ann Clin Microbiol Antimicrob* **16**:10.
- 434 34. **Neilson AA, Mayer CA.** 2010. Cholera - recommendations for prevention
435 in travellers. *Aust Fam Physician* **39**:220-6.
- 436

437 **Table and Figures**

438

439 Table 1. Summary of GeneFinder profiles and ST results

440

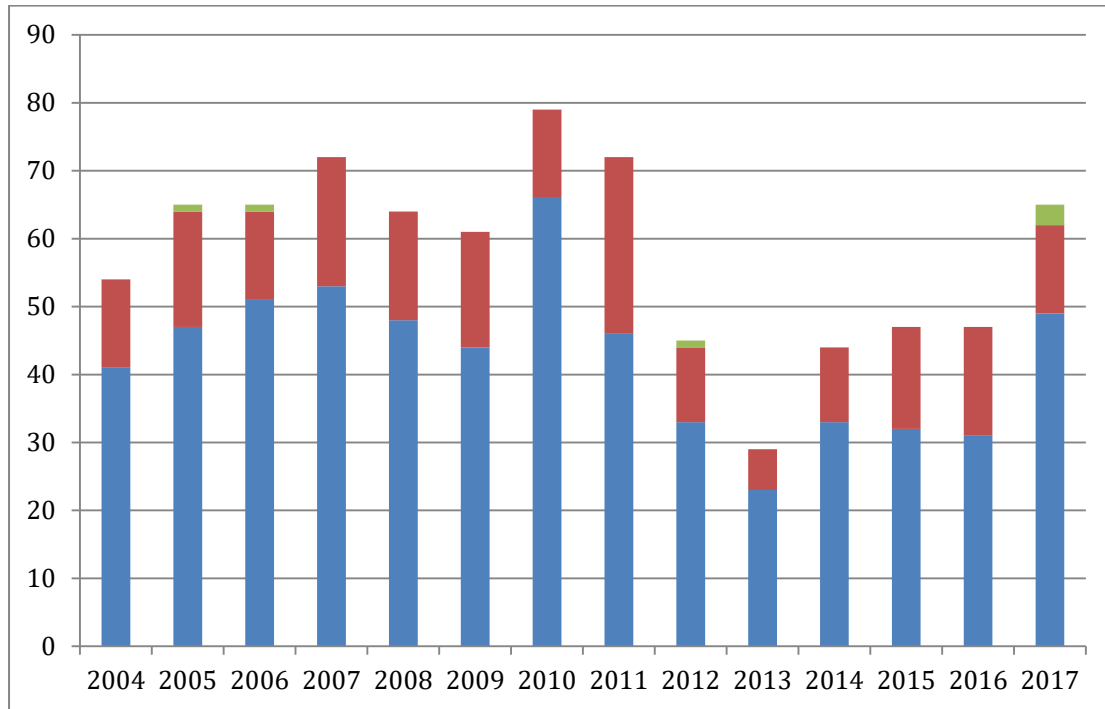
Genefinder profile	ST	Number of isolates
<i>toxR, wbeO1, tcp Classical, ctxA</i>	73	3
<i>toxR, wbeO1, tcp El Tor, ctxA</i>	69	34
<i>toxR, wbeO1, tcp El Tor</i>	75, 169, 579 (2)	4
<i>toxR, wbeO1, ctxA</i>	167	1
<i>toxR, wbeO1</i>	521(2), 551 (2)	4
<i>toxR, wbfO139, tcp El Tor, ctxA</i>	69	1
<i>toxR, wbfO139</i>	163, 527, 529, 544, 568	5
<i>toxR (77%), wbfO139</i>	586	1
<i>toxR</i>	>70 different STs	94
<i>toxR (77%)</i>	539, 540, 541, 550, 585, 587, 600	7
	Total	152

441

442

443

444 Figure 1. Number of isolates of *V. cholerae* from human cases resident in England submitted
445 to GBRU by local hospital laboratories each year between 2004 and 2017 (n=836). Non-O1,
446 non-O139 serogroups - blue; Serogroup O1 - red; Serogroup O139 - green
447



448

449

450 Figure 2. Minimum spanning tree illustrating the diversity in the population structure of the
451 isolates of *V. cholerae* received at PHE between 2015 and 2017. Clonal complexes (CC)
452 comprising strains linked by a single locus variant (thick black line) or double locus variant
453 (thin black line) and are shaded grey. Sequence types (ST) are shown in black. Isolates
454 associated with cases reporting recent travel abroad are highlighted: red – Asia; blue –
455 Africa, green –Latin America, yellow – mainland Europe, white – no data.

456

457

458 Supplementary Figure 1. Phylogeny of ST69 comprising isolates from this study (highlighted
459 in red) and from publicly available databases

460

461 Supplementary Table 1. Short read archive accessions, WGS data and travel data for the
462 sequenced isolates (n=152)

463

