1    **05 July 2022**

2    **Title:**

3    **Genomic epidemiology of the cholera outbreak in Yemen reveals the spread of a multi-**

4    **drug resistance plasmid between diverse lineages of *Vibrio cholerae***

5

6    Florent Lassalle[1]*, Salah Al-Shalali[2], Mukhtar Al-Hakimi[2], Elisabeth Njamkepo[3], Ismail

7    Mahat Bashir[4], Matthew J. Dorman[1], Jean Rauzier[3], Grace A. Blackwell[1,5], Alyce Taylor-

8    Brown[1], Mathew A. Beale[1], Ali Abdullah Al-Somainy[6], Anas Al-Mahbashi[2], Khaled

9    Almoayed[7], Mohammed Aldawla[8], Abdulelah Al-Harazi[6], Marie-Laure Quilici[3$], François-

10   Xavier Weill[3$], Ghulam Dhabaan[9]*$, Nicholas R. Thomson[1,10]*$

11

12   $ Joint last authors

13   *corresponding authors

14

15   *Affiliations*

16   [1] Parasites & Microbes Programme, Wellcome Sanger Institute, Wellcome Genome Campus,

17   Hinxton, Cambridge, CB10 1SA, United Kingdom

18   [2] Faculty of Science, Sana'a University, Sana'a, Yemen

19   [3] Institut Pasteur, Université Paris Cité, Unité des Bactéries pathogènes entériques, F-75015,

20   Paris, France

21   [4] WHO Yemen country office, Sana'a, Yemen

22   [5] EMBL-EBI, Wellcome Genome Campus, Hinxton, United Kingdom

23   [6] National Centre of Public Health Laboratories, Sana'a, Yemen

24   [7] Ministry of Public Health and Populations, Diseases Control & Surveillance, Sana'a, Yemen

25   [8] Ministry of Public Health and Populations, Infection Control Unit, Sana'a, Yemen

26   [9] Department of laboratory Medicine and Pathobiology, Faculty of Medicine, University of

27   Toronto, Canada

28   [10] London School of Hygiene and Tropical Medicine, London, WC1E 7HT, United Kingdom

29

30

31   **Abstract**

32

33   The humanitarian crisis in Yemen led in 2016 to the biggest cholera outbreak documented in

34   modern history, with more than 2.5 million suspected cases to date. In late 2018,

35   epidemiological surveillance showed that *V. cholerae* isolated from cholera patients had turned

36   multi-drug resistant (MDR). We generated genomes from 260 isolates sampled in Yemen

37   between 2018 and 2019 to identify a possible shift in circulating genotypes. 84% of *V. cholerae*

38   isolates were serogroup O1 belonging to the seventh pandemic El Tor (7PET) lineage,

39   sublineage T13 – same as in 2016 and 2017 – while the remaining 16% of strains were non-

40   toxigenic and belonged to divergent *V. cholerae* lineages, likely reflecting sporadic gut

41   colonisation by endemic strains. Phylogenomic analysis reveals a succession of T13 clones,

42   with 2019 dominated by a clone that carried an IncC-type plasmid harbouring an MDR pseudo-

43   compound transposon (PCT). Identical copies of these mobile elements were found

44   independently in several unrelated lineages, suggesting exchange and recombination between

45  endemic and epidemic strains. Treatment of severe cholera patients with macrolides in Yemen
46  from 2016 to early 2019 coincides with the emergence of the plasmid-carrying T13 clone. The
47  unprecedented success of this genotype where an SXT-family integrative and conjugative
48  element (SXT/ICE) and an IncC plasmid coinhabit show the stability of this MDR plasmid in
49  the 7PET background, which may durably reduce options for epidemic cholera case
50  management. We advocate a heightened genomic epidemiology surveillance of cholera to help
51  control the spread of this highly-transmissible, MDR clone.

52

53  **Introduction**

54

55  Since 2016, Yemen has seen the largest epidemic of cholera ever recorded. This occurred
56  against the backdrop of a civil war turned international conflict and famine which together
57  fueled extensive population movement, with more than 4 million people internally displaced
58  by the end of 2020[1]. The Electronic Disease Early Warning System (eDEWS), a surveillance
59  programme coordinated by the Ministry of Public Health and Population of Yemen (MPHP) in
60  Sana'a tasked with monitoring the epidemic[2], had recorded a total of almost 2.4 million
61  suspected cholera cases up until August 2019[3]. These cases exhibited a seasonal profile, with
62  peaks in July 2017 and September 2018 (16,000 and 50,000 cases per week, respectively)[3].
63  The lower reported case incidence in 2018 was ascribed to the mass vaccination campaign led
64  by the World Health Organization (WHO) and United Nation Children's Fund (UNICEF), who
65  delivered the oral cholera vaccine (OCV) to 540,000 people in August 2018 (387,000 at follow-
66  up in September) in targeted districts in Aden, Hudaydah and Ibb governorates[4,5].
67  Notwithstanding this focussed vaccination campaign, cholera cases were recorded nationwide
68  in 2019, peaking at over 30,000 cases per week. Despite the mass vaccination campaign, case
69  numbers declined at a slower rate than in previous years[3].

70

71  Pandemic cholera is caused by discrete phylogenetic lineages of the bacterium *Vibrio cholerae*
72  that are associated with epidemic spread, and carry lipopolysaccharide O-antigens of
73  serogroups O1 or O139. The large majority of epidemic strains associated with cholera
74  outbreaks from the last 60 years belong to the seventh pandemic El Tor (7PET) lineage of *V.*
75  *cholerae* O1, which swept the planet in three pandemic waves[6]. We previously used genomic
76  epidemiology to show that the first two waves of the cholera outbreak in Yemen (2016 and
77  2017) were driven by a single clonal expansion[7] belonging to Wave 3 of the global 7PET
78  lineage and had an Ogawa serotype. This indicated the Yemen outbreak was seeded by a single
79  international transmission event linked to the 7PET sublineage involved in the thirteenth
80  recorded intercontinental introduction of cholera (T13)[7].

81

82  Our ongoing surveillance activities in Yemen found that the fluctuating peaks in incidence in
83  Yemen were accompanied by a sudden change in the antibiotic susceptibility profile reported
84  by the reference laboratory at the MPHP in Sana'a. Whilst strains isolated in 2016-2018 were
85  sensitive to most of the antibiotics usually used for the treatment of cholera (excepting
86  quinolones, where reduced suceptibility to ciprofloxacin prevented the use of this antibiotic as
87  a single dose treatment), by 2019, resistance was observed for multiple drugs including third
88  generation cephalosporins, macrolides (including azithromycin) and cotrimoxazole. Whilst the

89   main treatment for cholera is rehydration therapy, antibiotics can be used to limit the volume
90   and duration of the acute watery diarrhoea, and reduce the risk of transmission[8–10]. In Yemen,
91   macrolides were used extensively up to early 2019 to treat moderate to severe cases of cholera
92   in pregnant woman and children, the latter forming the large majority of cases[11]. Multiple drug
93   resistance (MDR) in *V. cholerae* is strongly associated with the acquisition of mobile genetic
94   elements (MGEs) such as SXT-family integrative and conjugative elements (SXT/ICE) or
95   plasmids of the incompatibility type C (IncC, formerly known as IncA/C2; ref. 12), which often
96   carry and disseminate antimicrobial resistance (AMR) gene cargo[13].
97
98   We hypothesised that the MDR phenotype seen in the Yemen *V. cholerae* isolates from 2019
99   could be explained either by gain of resistance (either through *de novo* mutations or acquisition
100  of resistance-conferring MGEs) in the previously susceptible 7PET-T13 *V. cholerae* strain
101  already circulating in Yemen, or through the replacement of that strain with locally or globally
102  derived MDR strain(s). Distinguishing between these hypotheses is important for
103  understanding the ongoing dynamics of cholera in Yemen, and will be important for cholera
104  control strategies. We therefore applied genomic epidemiology approaches to determine the
105  molecular basis for the observed switch to the MDR phenotype and its link to global and local
106  evolutionary dynamics of pandemic cholera. In doing so, we highlight the role of globally
107  circulating MGEs in making an epidemic pathogen resistant to multiple drugs and subsequently
108  reducing treatment options. We also show that these MGEs and their cargo AMR genes were
109  repeatedly exchanged among diverse *V. cholerae* lineages found in Yemen.
110
111
112  **Results**
113
114  **Sampling of *V. cholerae* in Yemen in 2018 and 2019**
115
116  The National Centre of Public Health Laboratories (NCPHL) in Sana'a, the capital city,
117  received 6,311 and 3,225 clinical samples collected from suspected cholera patients, in 2018
118  and 2019 respectively. Of these, 2,204 (35%) and 2,171 (67%) were confirmed to be positive
119  for *V. cholerae* O1 by culture (identification based on biochemical tests and detection of Ogawa
120  and Inaba serotypes; Table S7; Figure S1). Among the 1,642 *V. cholerae* isolated at the NCPHL
121  from January to October 2018, 623 were tested for susceptibility to a range of antibiotics by
122  the disk diffusion method, of which 620 (99.6%) were phenotypically resistant to nalidixic acid
123  and nitrofurantoin, but otherwise sensitive to all other antimicrobials tested (Figure S2; Tables
124  S7). In contrast, all tested *V. cholerae* isolates (*n* = 2,172) from January 2019 onwards were
125  resistant to nalidixic acid, azithromycin, co-trimoxazole and cefotaxime (Figure S2; Tables
126  S7), a pattern maintained up to late 2021 (WHO EMRO, personal communication). The
127  transition in phenotype occurred during November 2018, when 159/175 (90.8%) tested isolates
128  already showed the MDR profile. 250 of the 2018-2019 clinical *V. cholerae* isolates were
129  randomly chosen for further characterization (Table S1). These samples originated from eight
130  of the 21 Yemen governorates, comprising 71 out of 333 districts (Table S1), with 101 samples
131  collected in 2018 (from mid-July to late October) and 149 in 2019 (from late February to late

132    April and from early August to mid-October). In addition, ten environmentally-derived strains
133    were isolated from sewerage in Sana'a in October 2019 (Table S1).

134

135    Extended antibiotic sensitivity testing of these 260 isolates at NCPHL and Institut Pasteur (IP)
136    (Figure S3) reflected the phenotypic switch to MDR observed in the wider sample set, further
137    showing that all tested 2019 strains were resistant to ampicillin, cefotaxime, nalidixic acid,
138    azithromycin, erythromycin and co-trimoxazole (Tables S1, S2; Supplementary text).
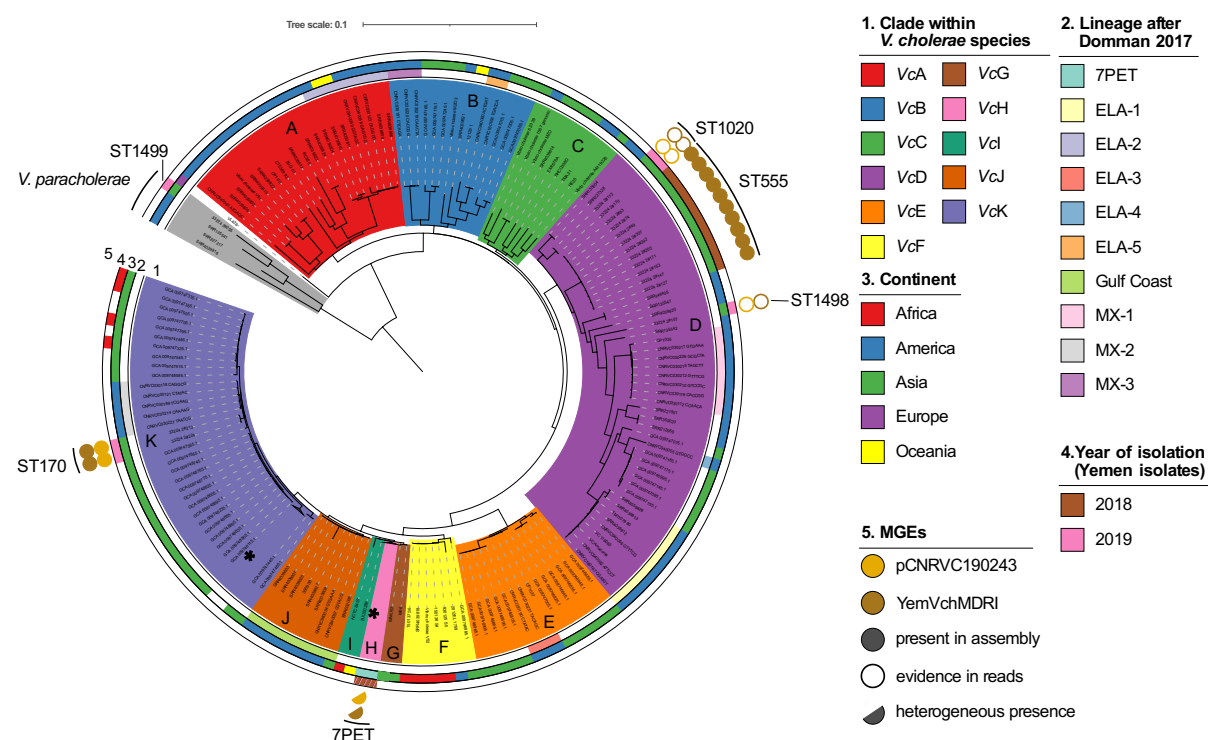
139

140    **Phylogenetic diversity of the *V. cholerae* isolated in Yemen in 2018 and 2019**

141

142    We isolated a single colony for 240 out of the 260 *V. cholerae* isolates indicated above, and
143    multiple independent colony picks for the remaining 20 isolates, for a total of 281 isolates on
144    which we performed whole genome sequencing (Figure S3; Tables S1, S2). After quality
145    filtering, this yielded 232 high-quality isolate genome assemblies (selecting a single isolate
146    from each initial sample), which we combined with 650 previously published *V. cholerae* O1
147    and non-O1 genomes for context, totaling 882 assembled genomes (Table S4; Figure S3). We
148    inferred a core-genome phylogeny for this genome set, which described the sequenced diversity
149    of the *V. cholerae* species, rooted by the genomes that belong to its newly described sister
150    species *V. paracholerae*[14]. We subdivided *V. cholerae* genomes according to their distribution
151    in eleven crown clades of the core-gene phylogeny clades, referred to henceforth as *Vc*A to
152    *Vc*K (Figures 1, S3; Table S4). *Vc*H contained all 7PET epidemic lineage genomes utilised in
153    this dataset, including 663 contextual genomes, the the majority (216/232) of the Yemen 2018-
154    2019 genomes, and all 42 previously reported 2016-2017 Yemeni genomes[7] (Figure S4).

155

156    Whilst Yemeni *Vc*H isolates show limited genomic diversity (99.98-100.00% ANI similarity;
157    0 to 97 SNPs), the remaining 16 Yemeni genomes belonged to clades *V. paracholerae* (*Vpc*),
158    *Vc*D and *Vc*K and were overall more diverse than *Vc*H isolate genomes (96.24-99.99% ANI
159    similarity; Figure 1; Table 1); these represent "non-7PET" lineages. Based on core genome
160    phylogeny and MLST, we found five distinct clusters within three non-7PET clades: *Vpc* ($n =$
161    1; novel ST1499), *Vc*D ($n = 21$; ST555, ST1020 and novel ST1498; Table S6) and *Vc*K ($n = 2$;
162    ST170) (Figure 1).

163

164    Although highly clonal, phylogenetic structure within *Vc*H allowed it to be further subdivided
165    into subclades *Vc*H.1 to *Vc*H.10 (Figure S5). All the Yemen 2016-2019 isolates fell within
166    *Vc*H.9, which corresponds to the T13 sublineage of 7PET Wave 3 (ref. 7). We selected one
167    representative isolate (CNRVC190243) of *Vc*H.9, and used PacBio sequencing to generate long
168    reads in addition to the Illumina short reads obtained for all samples, which enabled us to
169    generate a closed hybrid assembly. We subsequently used Oxford Nanopore sequencing to do
170    the same for a *Vc*D representative isolate (CNRVC190247). To obtain greater phylogenetic
171    resolution within *Vc*H.9, we then mapped sequencing reads to our new *Vc*H.9 CNRVC190243
172    reference genome to build a "mapped genome tree". Here, together with our novel *Vc*H.9
173    genomes ($n = 238$), we included 218 previously published genomes that reside in this subclade
174    and close outgroups, for a total of 456 genomes (Table S5). This approach allowed us to further
175    subdivide *Vc*H.9 into phylogenetic clusters named *Vc*H.9.a to *Vc*H.9.h (Figure 2A). Yemeni

genomes form a monophyletic group (clusters *Vc*H.9.e to *Vc*H.9.h), emerging from the genetic diversity of East African genomes (clusters *Vc*H.9.c and *Vc*H.9.d), which in turn branch out of a cluster of South Asian genomes (*Vc*H.9.b), consistent with previous observations on the origins of 7PET-T13, introduced from South Asia into Africa [7,15]. Clusters *Vc*H.9.g and *Vc*H.9.h together comprise the majority of 2018-2019 Yemen isolates (235/281) and form a well-supported clade (94% bootstrap) that branches from within *Vc*H.9.f (Table 1). Cluster *Vc*H.9.h includes the majority of the Yemeni 7PET-T13 isolates (78/87) from 2018, with just one isolate from March 2019. In contrast, Cluster *Vc*H.9.g comprises mostly 2019 isolates (150/156), and a minority from 2018 (6/156) (Table 1). All of the 2016-2017 Yemen isolates (*n* = 42) belong to sister clusters *Vc*H.9.e and *Vc*H.9.f.



**Figure 1: Phylogenetic diversity of *Vibrio cholerae* isolates from Yemen**
Maximum-likelihood phylogeny of 882 assembled *V. cholerae* genomes based on the 37,170 SNP sites from the concatenated alignments of 291 core genes. Low-diversity clades (*Vc*H and part of *Vc*K) are collapsed and marked by black stars. Clades are highlighted with background colours (legend key 1). Coloured rings outside the tree depict the match with previously described lineages (ring 2), the geographical origin of isolates at the level of continents (ring 3), and their year of isolation when from Yemen (ring 4). Presence of parts of the plasmid pCNRVC190243 are indicated by coloured circles (ring 5 in A): IncC plasmid backbone (light brown) and the MDR pseudo-compound transposon Yem*Vch*MDRI (dark brown); full circles indicate over 70% coverage in assemblies of the reference length, hollow circles indicate 30-70% coverage in assemblies and confirmed presence based on mapped reads, with even coverage over the MGE reference sequence, while half-circles represent heterogeneous presence in a collapsed clade. Tree plots were generated with iTOL v4[16] and adapted with Inkscape. The scale bar represents the number of nucleotide substitutions per site.

**Spatiotemporal distribution of *V. cholerae* isolates**
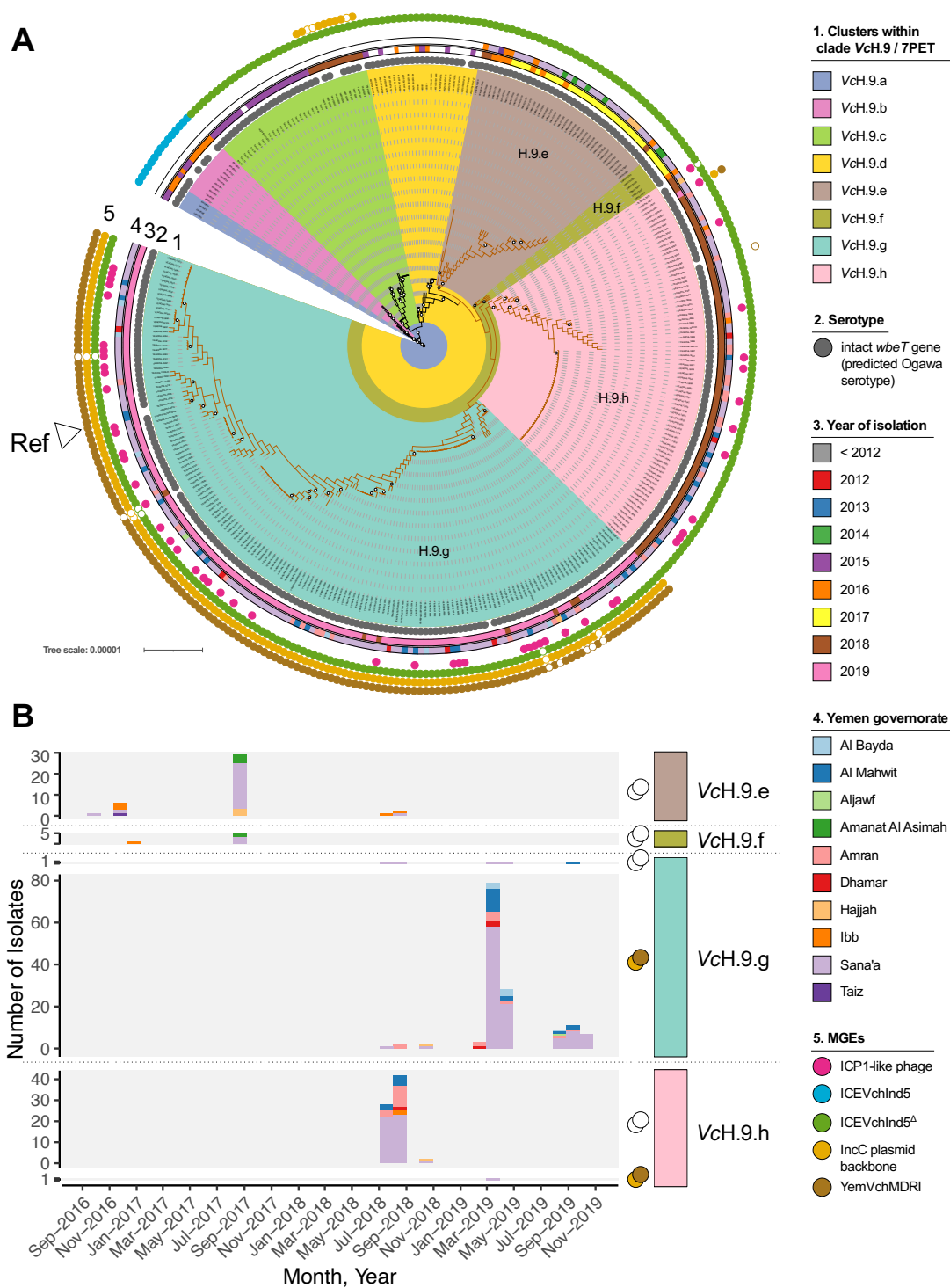
To delineate the evolutionary dynamics of the cholera outbreak in Yemen, we plotted *Vc*H.9 isolates by phylogenetic cluster over time (based on the date of sample collection) and between

205  administrative divisions (linked to reporting hospital). From Figure 2B it is clear that each
206  annual wave was dominated by a single cluster: 2016 and 2017 by *Vc*H.9.e; 2018 by *Vc*H.9.h;
207  2019 by *Vc*H.9.g. There was no evidence of geographic restriction for any of these clusters,
208  even when accounting for dispersal over time (Fig 2C, 2D; Table S5; Supplemental data online,
209  doi: 10.6084/m9.figshare.19097111). Next, we analysed the relationship between temporal and
210  spatial distances, based on the date and GPS coordinates of sample collection, as well as with
211  the pairwise phylogenetic distances between genomes. We found no significant correlation
212  between the spatial and temporal distances, nor between the spatial and phylogenetic distances
213  (Table S8). These data did show a positive correlation between the temporal and phylogenetic
214  distances ($R^2$ = 0.181; Mantel test *p*-value < $10^{-6}$) (Table S8), with root-to-tip distances
215  significantly correlated with sampling date (Pearson's $R^2$ = 0.437; $p < 10^{-15}$).
216
217  **Figure 2: Phylogenetic diversity and spatiotemporal distribution of *Vibrio cholerae* 7PET-T13 isolates**
218  **(*Vc*H.9) from Yemen**
219  **A.** Subtree of the maximum-likelihood phylogeny of 456 7PET genomes mapped to reference *Vc*H.9 strain
220  CNRVC190243 genome, including 335/456 genomes covering *Vc*H.9 (as defined in Figure S5), which
221  corresponds to the 7PET-T13 sublineage and close South Asian relatives. The full tree containing the 456 genomes
222  is available as supplementary material on Figshare (https://figshare.com/s/4d83a32cce78a52b413e; doi:
223  10.6084/m9.figshare.16595999) and was obtained based on 2,092 SNP sites from concatenated whole-
224  chromosome alignments. Brown branches indicate the clade grouping all Yemeni 7PET-T13 isolates. Bootstrap
225  support over 70% is indicated by white circles. Phylogenetic clusters within *Vc*H.9 are highlighted with
226  background colours (legend key 1). Coded tracks outside the tree depict the serotype of isolates (ring 2) as
227  predicted from genomic data, year of isolation when isolated in 2012 or later (ring 3), the governorate of isolation
228  if in Yemen (ring 4). The presence of mobile genetic elements (MGEs) is indicated by coloured circles in the
229  outermost track (ring 5): ICP1-like phage (pink), SXT/ICE ICE*Vch*Ind5 (blue), ICE*Vch*Ind5$^\Delta$ i.e. featuring the
230  characteristic 10-kb deletion in the variable region III (green), IncC plasmid backbone (light brown) and the MDR
231  pseudo-compound transposon Yem*Vch*MDRI (dark brown); filled and unfilled circles indicate different level of
232  coverage in assemblies (see Figure 1 legend). The position of the reference sequence to which all other genomes
233  were mapped to generate the alignment is labelled. The scale bar represents the number of nucleotide substitutions
234  per site. **B.** Frequency of each phylogenetic subcluster among Yemen isolates per month since the onset of the
235  Yemen outbreak. Where relevant, the cluster group is subdivided by the presence or absence of the IncC plasmid
236  as indicated by the filled brown (present) or open (absence) circle on the right of the chart. The contribution of
237  each governorate of isolation is indicated by the coloured portion of each bar. **C** and **D.** A map of Yemen
238  governorates (C) and a focus on the Sana'a and Amanat Al Asimah governorates (inner and outer capital city; D),
239  with dots corresponding to isolates, coloured by phylogenetic subcluster.
240

241

242

243 We inferred a timed phylogeny for *Vc*H.9 (Figure S6), which revealed that the most recent

244 common ancestor (MRCA) of all Yemeni *V. cholerae* 7PET-T13 genomes was estimated to

245 have existed in February 2015 (95% confidence interval [95%CI], April 2014 and July 2015).

246 Moreover, the MRCAs for clusters *Vc*H.9.e and *Vc*H.9.f (mostly sampled in 2016 and 2017)

247 were dated May and June 2015, respectively, and the MRCAs for clusters *Vc*H.9.g and *Vc*H.9.h

248 (sampled in 2018 and 2019) were dated February and March 2017, respectively. In addition,

249 we dated the MRCA of the clade grouping clusters *Vc*H.9.g and *Vc*H.9.h, which represent the

250 majority of 2018-2019 Yemen isolates, to September 2016 (Figure S6).

251

252 The distribution of non-7PET isolates across Yemen was mostly sporadic (Table S4;

253 Suplementary Text; Supplementary data online https://figshare.com/s/73fcd5e1b4958c97ef78,

254 doi: 10.6084/m9.figshare.19097111). However, we characterised a cluster of eighteen closely

255 related *Vc*D isolates belonging to ST555 (Table S6), which we found to differ from each other

256 by 0 to 10 SNPs (average 99.98% ANI similarity). Of these 18 isolates, 13 were isolated over

257 a period of 11 days in late July/early August 2018, two at the end of August, two in October

258 2018, and one in March 2019 (Table S6). They were obtained from patients in the neighbouring

259 governorates of Sana'a ($n = 7$), Al Mahwit ($n = 4$) and Amran ($n = 1$), which surround the

260 capital city. Genomes from other ST555 isolates, including strains reported as linked to

261 travelers returning to the UK from India in September 2015 and July 2016 (strains 229152 and

262 338360) [17], as well as closest relatives from our core-genome tree, were gathered to build a

263 mapped genome tree of *Vc*D genomes using the complete genome of 2018 Yemen strain

264 CNRVC190247 as a reference (Figure S7). The closest relative to Yemeni ST555 isolates,

265 strain 338360, differs from the *Vc*D ST555 genomes sequenced here by between 763-800

266 SNPs, ruling out direct clonal relationships.

267

268 **Predicted phenotypic properties of *V. cholerae* isolates**

269

270 Consistent with our previous report[7], Yemeni *Vc*H.9 isolates – which all belong to 7PET-T13

271 sublineage – all carried genes or mutations known to confer resistance to trimethoprim (*dfrA1*)

272 and to nalidixic acid (*gyrA*_S83I and *parC*_S85L). They also carried the *Vibrio* pathogenicity

273 island 1 (VPI-1, encoding the toxin co-regulated pilus TCP), VPI-2, the *Vibrio* seventh

274 pandemic islands I and II (VSP-I and VSP-II), and the CTX prophage, which all featured the

275 cholera toxin genes, *ctxAB,* of the allelic type *ctxB7*. None of the non-7PET genomes from

276 Yemen possessed a CTX prophage or the *ctxAB* genes. However, Yemni isolates belonging to

277 *Vc*K (ST170, related to previously described lineage MX-2), which were derived from the stool

278 of patients presenting cholera-like disease, carried all the genes coding for the TCP.

279

280 These *ctxAB⁻*, *tcpA⁺* *Vc*K genomes also carried the O1 LPS O-antigen biosynthetic gene cluster,

281 consistent with what has been seen previously in related non-7PET isolates[18]. The genomes of

282 the *Vc*D isolates belonging to ST555, ST1020, ST1498 and the *V. paracholerae* isolate

283 (ST1499), carried LPS O-antigen biosynthetic gene clusters encoding unknown serogroups;

284 these were conserved within and specific to each ST (Table S4; Figure S7A). Of the 216

285 Yemeni 2018-2019 *Vc*H isolates, 213 were predicted to produce a serogroup O1 LPS O-antigen

8

286 based on presence of a full biosynthetic gene cluster; in the three remaining assemblies this
287 genomic region was interrupted (YE-NCPHL-19012) or completely missing (YE-NCPHL-
288 18033 and YE-NCPHL-19140), likely due to limited genome sequence coverage (Table S3).
289 All predicted O1 serogroup isolates were predicted to be Ogawa serotype except two that
290 showed a disruption in *wbeT,* indicative of an Inaba phenotype (YE-NCPHL-18053 and YE-
291 NCPHL-19014, with gene truncation and point mutation respectively; Table S5). These
292 predictions were imperfectly reflected by the results of serological assays conducted at NCPHL
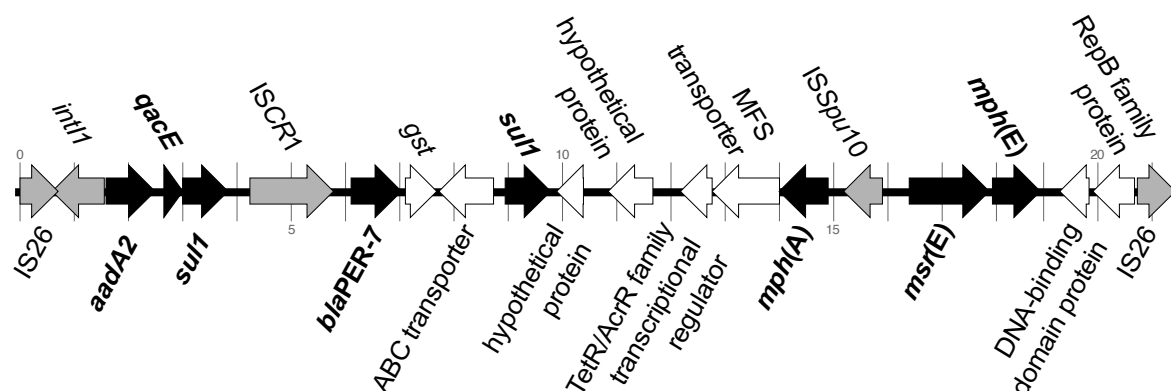293 (Table S2; Figure S8), suggesting issues in initial laboratory testing (see Supplementary Text).
294
295 **Genome variation of *Vc*H.9 (7PET-T13) isolates circulating in Yemen**
296
297 Given the change in antimicrobial suceptibility seen in the 2018-2019 Yemen isolates, we
298 compared in detail all of the *Vc*H.9 isolate genomes from Yemen to each other and related
299 isolates taken elsewhere, focusing on genotypic traits that were conserved in pandemic
300 sublineages occurring in Yemen. We identified three, four, and 21 fixed SNPs in the the crown
301 clade containing *Vc*H.9.e,f,g,h, the clade containing *Vc*H.9.g,h, and *Vc*H.9.h, respectively
302 (including 2, 2 and 11 non-synonymous SNPs, respectively; Table S9). Changes fell largely
303 within genes predicted to be involved in carbohydrate metabolism, signal transduction and
304 chemotaxis, none of which could be directly linked to change in virulence (Table S9).
305
306 Previously, the 2016-2017 Yemeni isolates carried an SXT ICE differing by only three or four
307 SNPs from the ICE*Vch*Ind5/ICE*Vch*Ban5 reference sequence (Genbank accession
308 GQ463142.1)[19], but which possessed a 10-kb deletion in variable region III, which explained
309 the phenotypic loss of resistance to streptomycin, chloramphenicol and sulphonamides (only
310 retaining resistance to trimethoprim via the *dfrA1* gene)[7]. All 2018-2019 *Vc*H.9 genomes
311 carried the same SXT ICE variant, with a maximum of 2 SNP differences and displaying the
312 same deletion. Hence, the change in antimicrobial resistance profile was not linked to variation
313 in SXT ICE.
314
315 Looking across all genes within the pangenome, the only variation directly associated with the
316 Yemen 2018-2019 genomes, compared to those sequenced from 2016-2017, was the presence
317 of a novel 139-kb plasmid, which we named pCNRVC190243 (Table S10). The backbone of
318 this new plasmid includes a replicon of the IncC type, as well as genes encoding a complete
319 type F conjugative apparatus and a MOBH-type relaxase, suggesting it is self-transmissible.
320 Plasmid pCNRVC190243 also carries a 20-kb genomic region (which we denoted
321 Yem*Vch*MDRI) predicted to encode a quaternary ammonium compound efflux pump (*qac*), an
322 extended-spectrum beta-lactamase (ESBL; *bla*PER-7), sulphonamide resistance (*sul1*),
323 aminoglycoside resistance (*aadA2*), and macrolide resistance (*mph*(A), *mph*(E) and *msr*(E))
324 (Figure 3; Table S4). Yem*Vch*MDRI is a pseudo-compound transposon (PCT) – a structure
325 bounded by IS*26* elements[20] – and includes a class 1 integron with *aadA2* encoding resistance
326 to streptomycin and spectinomycin as a gene cassette, associated with an IS*CR*1 element
327 carrying the ESBL *bla*PER-7 gene, a structure similar to one previously seen in *Acinetobacter*
328 *baumanii*[21,22]. We found that pCNRVC190243 was present in 6/89 (6.7%) Yemeni *Vc*H.9
329 isolates from 2018, but this rose to 100% (151/151) in 2019 (Figure 2B). This was linked to

330  phylogenetic cluster, with only 1/79 (1.3%) *Vc*H.9.h isolates harbouring the plasmid, compared
331  to all (156/156) *Vc*H.9.g isolates (Figure 2A).

332



333
334  **Figure 3: Genetic organisation of the MDR pseudo-compound transposon Yem*Vch*MDRI**
335  Antimicrobial resistance (AMR) genes are filled in black and labelled in boldface; genes encoding endonucleases
336  transposases and other genes involved in genetic mobility are filled in grey. Genomic position is indicated by
337  tickmarcks every kilobase, in reference to the pCNRVC190243 plasmid coordinates.

338

### Distribution and relatedness of MDR mobile genetic elements

340  Analysis of the broader phylogenetic context of pCNRVC190243 suggested at least three
341  independent acquisitions of this plasmid (and associated Yem*Vch*MDRI), since it was also
342  present in three *Vc*D (ST1499 and ST1020) and two *Vc*K (ST170) isolates collected in 2019 in
343  Yemen. Comparing the full-length sequence of all pCNRVC190243 plasmids from *Vc*H.9 and
344  *Vc*K and *Vc*D isolates showed all sequences were identical except for two isolates: one varied
345  by a single SNP resulting in an amino-acid change S71F in the sulphonamide resistance protein
346  Sul1 (YE-NCPHL-19105; G26720A SNP); the other by a single intergenic SNP. We also found
347  the Yem*Vch*MDRI element integrated into chromosome 2, without the pCNRVC190243
348  backbone (Figure 1; Supplementary Text) in all eighteen of the ST555 isolates.

349  Searching a broader prokaryotic genome database, closely related but non-identical elements
350  were found in different combinations in other *V. cholerae* and diverse bacterial taxa: an IncC
351  plasmid, named pYA00120881 (GenBank accession MT151380), was identified in 13 closely
352  related *Vc*H.9.a and *Vc*H.9.c isolates (Figure 2A) that were collected in 2015 and 2018 in
353  Zimbabwe[15]. The backbones of these IncC plasmids share 99.98% nucleotide sequence
354  identity, but pYA00120881 carries a different MDR genomic region – featuring a *bla* gene
355  encoding a CTX-M-15 ESBL – inserted at the same locus (Figure S9). Furthermore, 59 *V.*
356  *cholerae* O139 (ST69) isolates collected in China from 1998 to 2009[23,24] (unpublished genomic
357  data released in BioProject PRJNA303115; Table S11) carry IncC-type plasmids that show
358  similarity to pCNRVC190243 and also include Yem*Vch*MDRI-like PCT elements, albeit
359  lacking IS*CR*1 and its associated *bla*$_{PER-7}$ gene.

360

361  Importantly, when using the Yem*Vch*MDRI sequence alone for database searches, we found
362  the genome of *V. cholerae* ST555 strain 338360 (Table S6) shared 100% nucleotide identity

10

363    with the complete Yemeni ST555 Yem$Vch$MDRI sequence, including the $bla_{PER-7}$-carrying
364    IS$CR$1 (IS$CR$1$_{blaPER-7}$; Table S12). Likewise, IS$CR$1$_{blaPER-7}$ has also been previously observed
365    in the genomes of *A. baumanii* strains[21,25] from France and the United Arab Emirates (UAE).
366    Those from UAE were located on the plasmid pAB154, where the sequence homology with
367    IS$CR$1$_{blaPER-7}$ extended beyond the canonical element and included Yem$Vch$MDRI flanking
368    regions, suggesting that the IS$CR$1$_{blaPER-7}$ carried by pAB154 is derived from Yem$Vch$MDRI,
369    or a closely related element (Figure S10). Moreover, outside of *V. cholerae,* pCNRVC190243-
370    and/or Yem$Vch$MDRI-like elements are widely distributed with *Escherichia coli, Salmonella*
371    *enterica* and *Klebsiella pneumoniae* genomes presenting >95% shared nucleotide *k*-mers
372    (Table S11, S12; Figure S11), with the closest matches outside of *V. cholerae* being seen in *K.*
373    *pneumoniae*. This indicates that similar regions may be widely distributed in MGEs across
374    bacterial taxa.
375
376    A recent study reported that two anti-plasmid defence systems, DdmABC and DdmDE, cause
377    the instability of plasmids in *V. cholerae* cells, including those with IncC-type replicons[26].
378    These proteins are encoded by all 7PET genomes, which according to this study would explain
379    the incapacity of plasmids to be maintained in a 7PET background. We here show that a 139-
380    kb IncC-type plasmid has been stably propagated in a clone of the 7PET lineage. We verified
381    the presence and integrity of the DdmABC and DdmDE systems and found they were present
382    and intact in all 7PET genomes in our 882 assembled genome dataset – including those
383    harbouring pCNRVC190243 (Table S5). This shows that these defence systems are not
384    sufficient to destabilise pCNRVC190243 to the point of it being lost from the population within
385    the 15-month period covered by our study, or even for the following two years, as suggested
386    by 2021 antibiotic susceptibility profiles.
387
388
389    **Discussion**
390
391    Characterising the genomic nature of the pathogens causing an outbreak can reveal changing
392    epidemiological dynamics through adaptive evolution of the pathogen or the introduction of
393    distinct pathogen lineages. Such events may lead to the emergence of a more virulent or drug
394    resistant genotype of the pathogen, and impact disease control efforts. Our genomic
395    epidemiology analysis shows that despite seasonal fluctuation, the vast majority of cholera in
396    Yemen is caused by the 7PET-T13 lineage (*Vc*H.9), and is derived from a single introduction
397    into Yemen. Using a larger sample set, we refined our previous date estimate[7] and show that
398    the progenitor of the Yemeni outbreak emerged between April 2014 and September 2015,
399    contemporarily to the onset of the civil war in Yemen, and existed one to two years prior to the
400    declaration of a cholera outbreak[3,7].
401
402    Using our high-resolution phylogenomic tree, we were able to subtype the majority of Yemeni
403    genomes into four different phylogenetic clusters that dominated the outbreak at different
404    points in time. We observed two large clonal expansions for the sister clades that dominated
405    2018 and 2019 (clusters *Vc*H.9.h and *Vc*H.9.g, respectively), which both emerged in early
406    2017. Founding effects at the onset of each cholera season, associated with rapid expansions,

11

407  may explain the dominance of each cluster in these respective epidemic waves. However, it is
408  also possible that an adaptive advantage participated in driving the replacement of *Vc*H.9.h by
409  its sister clone, *Vc*H.9.g. In the absence of samples from subsequent years (with surveillance
410  efforts hampered by the Covid-19 pandemic), it was not possible to establish whether these or
411  another MDR lineage persisted past 2019.

413  In Yemen, pregnant women and children (one third of cholera patients were aged 15 or under;
414  Table S7; ref. 11) were treated with erythromycin and azithromycin between 2016 and late
415  2018. The 2019 wave of the Yemen cholera outbreak was associated with a sudden change in
416  antibiotic resistance profile, from being largely sensitive to antimicrobials between 2016-2018,
417  to being resistant to multiple therapeutically relevant drugs in 2019. Our data showed that this
418  phenomenon coincided with the appearance in late 2018 of plasmid pCNRVC190243 in
419  isolates belonging to *Vc*H.9.g, the phylogenetic cluster which dominated our 2019 samples.
420  Plasmid pCNRVC190243 carries the pseudo-compound transposon Yem*Vch*MDRI, which in
421  turn includes a type 1 integron and the IS*CR*1$_{bla\text{PER-7}}$ element. These elements confer resistance
422  to third-generation cephalosporins, aminoglycosides, macrolides and sulphonamides (and,
423  combined with the *dfrA1* gene present on the SXT ICE, to co-trimoxazole), plus disinfectant
424  tolerance provided by the *qac* gene[27]. The acquisition of Yem*Vch*MDRI element by an ancestor
425  of *Vc*H.9.g was followed by its dramatic spread – a clonal expansion which we show to occur
426  in 2018 (Figure S6), a time when there would have been a selective pressure towards macrolide
427  resistance in symptomatic cases due to the large-scale administration of these drugs.

429  We also identified a small number of non-7PET *V. cholerae* amongst 2018-2019 Yemen
430  isolates: 8% of unique clinical isolates (21/254) and 30% of environmental isolates (3/10)
431  belonged to three diverse lineages. The location and times of isolation of these non-7PET *V.*
432  *cholerae* isolates suggest they largely represented sporadic infection events linked to endemic
433  strains. The only sizable cluster of non-7PET isolates were the 18 *Vc*D/ST555 isolates which
434  had near-identical genomes and isolated in 2018 and 2019 in several districts near the capital
435  city of Sana'a (Figure 1; Figure S7; Table S4). The short time range in which 15 of these strains
436  were isolated (31 days in July-August 2018) could be explained by repeated acquisitions from
437  a point source, although we cannot rule out that they stem from small-scale outbreaks, as has
438  been reported previously for non-O1/non-O139 strains[28,29]. However, we found no evidence of
439  long range spread of these non-pandemic clones across Yemen, characteristic of 7PET *V.*
440  *cholerae* isolates linked to epidemic disease. Importantly, the reappearance of this *Vc*D/ST555
441  genotype later in October 2018 and March 2019, with as little as 2 SNPs difference from
442  summer 2018 isolates, could suggest this genotype is able to persist in the environment,
443  possibly through similar ecological mechanisms as those that lead to the seasonal dynamics of
444  epidemic cholera following its initial introduction[11]. These ST555 strains might in fact
445  represent an endemic population of *V. cholerae* that can be carried without causing any disease.
446  These ST555 strains could have been isolated from a gut co-colonised by a cholera-causing
447  7PET strain, in an epidemic context where the pathogen is routinely isolated from cholera
448  patients using culture and enrichment techniques that are selective of the whole *V. cholerae*
449  species. This hypothesis of incidental isolation of ST555 strains in samples also containing a
450  toxigenic 7PET strain is supported by the original serotyping of all samples as O1 (Table S1;

451 Supplementary Text), and the positive detection of the *rfbO1* marker by PCR in samples from
452 which the four ST555 strains were isolated at the Institut Pasteur (IP) – samples which, when
453 sequenced at the Wellcome Sanger Institute (WSI), yielded 7PET genomes (Figure S3; Table
454 S13; Supplementary Text).
455
456 Sequences completely or near identical to plasmid pCNRVC190243, carrying the PCT
457 Yem*Vch*MDRI, were present in i) 7PET-T13 (*Vc*H.9) isolates, ii) all isolates from two of the
458 three different STs of *Vc*D (ST1499 and ST1020), and iii) the two *Vc*K/ST170 isolates; all of
459 which were collected in 2019 in Yemen. The only *Vc*H.9.h isolate from 2019 also carried this
460 plasmid. It is possible to explain these observations by multiple acquisitions of the plasmid
461 from independent sources or, more parsimoniously, as direct horizontal gene transfer events
462 between the *V. cholerae* lineages we report here. The large population sizes attained by the
463 epidemic lineages in Yemen make the latter hypothesis more likely, in a scenario of spill over
464 from the dominant cluster at the time, *Vc*H.9.g. However, we could not infer directionality due
465 to the limited available sampling from the diverse lineages in this study.
466
467 The Yem*Vch*MDRI PCT was also carried chromosomally by Yemeni ST555 strains. This PCT
468 is itself a composite element, including IS*CR1*$_{bla\text{PER-7}}$, a rare element that has only been
469 observed once in another *V. cholerae* background – the relatively closely-related ST555 strain
470 338360 (436 SNPs vs. reference strain CNRVC019247), isolated from a traveler returning from
471 India – and in two plasmids associated with *A. baumanii* strains isolated in Gulf countries.
472 Comparison of these homologous IS*CR1* elements suggests they all are derived from the same
473 ancestral element (Figure S10). Presence of the full, identical PCT Yem*Vch*MDRI in two
474 closely related, but distinct ST555 strains isolated from completely different geographical
475 origins, suggests this element is stably associated with this genotype. From this, we can
476 speculate that Yem*Vch*MDRI was originally present in Yemen in a ST555 genomic
477 background, and later combined with an IncC plasmid backbone to produce pCNRVC190243.
478 Again, directionality cannot be confidently inferred because of uneven sampling of host
479 lineages, and the potentially large number of unobserved donor bacteria.
480
481 Whilst pCNRVC190243 is a novel element, plasmids such as pYAM00120881 identified in
482 *Vc*H.9 *V. cholerae* from Zimbabwe in 2015 and 2018[15] shared almost identical plasmid
483 backbones. In addition, similar plasmids, some of which also carry Yem*Vch*MDRI-related
484 elements, have been observed in *V. cholerae* O139 isolates from China as well as detected in a
485 range of other bacterial genera, illustrating how widely distributed these IncC plasmids are.
486 Similarly, Yem*Vch*MDRI may occur in more diverse and more widely spread genomic
487 backgrounds that haven't been sampled yet. It is therefore possible that parts of this plasmid
488 have combined outside of Yemen from identical or similar genomic sources, independently
489 from the Yem*Vch*MDRI-carrying Yemeni ST555 strain.
490
491 While other MDR IncC plasmids were previously observed in *V. cholerae* in DRC, Kenya (in
492 a T10 sublineage genetic background) and Zimbabwe (T11 background in 2015), these were
493 only linked to sporadic cases or small-scale cholera outbreaks[30], despite selective conditions
494 linked to the widespread and uncontrolled use of antibiotics. Recently, it has been shown that

13

495    two defence systems, called DdmABC and DdmDE had the capacity to destabilise plasmids,
496    including large IncC-type plasmids[26]. These defence systems, encoded in all 7PET *V. cholerae*
497    genomes, were proposed to be responsible for the lack of maintenance of MDR plasmids in
498    populations of this pandemic lineage when not under stringent selective pressure for antibiotic
499    resistance. A first exception to the pattern of plasmid instability in 7PET *V. cholerae* was the
500    Zimbabwean cholera outbreak of 2018, which lasted six months and produced over 10,000
501    suspected cases, and was associated with a strain of the T13 background carrying the MDR
502    IncC plasmid pYA00120881[15]. The Yemeni cholera outbreak provides a further example, with
503    the T13 strain of the *Vc*H.9.g clone, carrying pCNRVC190243, being presumably associated
504    with more than a million suspected cases recorded since its emergence in late 2018. However,
505    the intact presence of the genes encoding the Ddm proteins in Yemeni and Zimbabwean
506    *Vc*H.9/7PET-T13 genomes presenting MDR IncC plasmids indicates there may be other
507    mechanisms that impact plasmid stability in 7PET genomes.
508

509    One possibility would be that an unknown environmental factor has applied a consistently
510    strong selective pressure for a trait carried by these plasmids. Even though the treatment of
511    cholera patients with macrolides was stopped in Yemen in early 2019, antibiotic pressure
512    remains a potential selective factor, as antibiotics and particularly azithromycin have been
513    reported to be overused by the general population in Yemen during the Covid-19 crisis[31].
514    Another possible factor would be the interaction with other mobile elements, including ICP1
515    phages, which we detected in a significant fraction of the samples (see Supplementary Text).
516    It has also been proposed previously that the presence of an SXT ICE in these genomes could
517    prevent the stable replication of an IncC-type plasmid, through an unknown functional
518    interference mechanism[7,15]. The unique occurrence of a 10-kb deletion in the SXT ICE
519    (ICE*Vch*Ind5/ICE*Vch*Ban5) in T13 isolates may provide these genomes with the novel capacity
520    to stably host an IncC plasmid; molecular genetic investigation of this locus should be
521    conducted to test whether it encodes another plasmid destabilisation factor. Whatever the
522    mechanism, it appears that both MDR elements SXT ICE and IncC plasmid are stably
523    propagated together in the Yemeni T13 strain, which population in Yemen has reached a
524    unprecedented size. This emerging MDR strain has therefore a high potential to spread and
525    seed further adapted lineages, as well as to disseminate its MDR plasmid and PCT to other
526    organisms.
527

528    **Conclusion**
529

530    The emergence of this multi-drug resistant pathogen demonstrates the necessity of continued
531    genomic surveillance of the microbial population associated with the ongoing Yemen cholera
532    outbreak, and for new outbreaks that may take place in regionally connected areas. Such
533    surveillance will enable Yemeni public health authorities to rapidly adapt clinical practices to
534    minimize AMR selective pressures. This also warrants increased efforts in research on the
535    molecular mechanisms and evolution of interactions between mobile genetic elements, to learn
536    about the constraints ruling their colonization of bacterial genomes. Such knowledge is
537    essential for us to be able to disentangle the role of MGEs from that of their bacterial hosts in
538    driving epidemics, so to propose practical definitions of pathogens that focus on the relevant

539 genes, mobile elements or prokaryotic organisms, and to implement appropriate molecular
540 epidemiology surveillance schemes.

541
542

543 **Materials and Methods**

544

545 **Definitions and surveillance data**
546 Cholera cases were notified to the the Ministry of Public Health and Population of Yemen
547 (MPHP) and recoded through the Electronic Disease Early Warning System (eDEWS)[2].
548 Suspected and confirmed cholera cases were defined according to the WHO in a declared
549 outbreak setting. Briefly, a suspected case is any person presenting with or dying from acute
550 watery diarrhoea (AWD) and a confirmed case is a suspected case with *Vibrio cholerae* O1 or
551 O139 infection confirmed by culture.

552

553 **Sample collection, microbiological testing and clinical metadata**
554 Clinical samples, i.e. stool and rectal swabs, were collected in Yemen by epidemiological
555 surveillance teams from suspected cholera cases during 2018 and 2019[11] and were transported
556 to the National Centre of Public Health Laboratories (NCPHL) in the capital city Sana'a in
557 Cary-Blair transport medium (Oxoid, USA). To probe the diversity of vibrios shed by
558 unreported cholera cases, as well as *V. cholerae* that may naturally occur in effluent waters,
559 environmental samples were collected during the day time in October 2019 from the sewage
560 system around Sana'a city and the vicinity and then transported to NCPHL for testing; each
561 sample was collected in sterile bottles containing enrichment media comprised of 250 mL of
562 sewage and alkaline peptone broth (APB, Difco Laboratories, Detroit, Michigan) at a 1:1 ratio
563 and incubated for 20 h at room temperature including the transportation time into the NCPHL
564 and processed as described previously[32]. All samples were cultured and identified according to
565 the Centers for Disease Control and Prevention (CDC) guidelines[33]. Resistance to antibiotics
566 was tested by the disk diffusion method according to the CLSI guidelines[34] for a range of
567 antibiotics as described in Table S1.
568 Live clinical isolates (n=120) were sent to the Institut Pasteur (IP; Paris, France), where only
569 21 samples were culture positive, due to poor sample preservation during shipment (Table S2;
570 Figure S3), leading to the final isolation of 22 *V. cholerae* strains (including two from mixed
571 culture YE-NCPHL-18020). Strains re-isolated at IP were characterized by biochemical and
572 serotyping methods according to standard practice of the French National Reference Centre for
573 Vibrios and Cholera (CNRVC)[35]. Separate antibiotic susceptibility testing (Table S2) was
574 performed by the disk diffusion method according to EUCAST guidelines (EUCAST 2020[36])
575 and MIC determination using the Sensititre™ (Thermo Scientific) and the Etest® (bioMérieux,
576 Marcy-l'Étoile, France) systems. Interpretation into S (Susceptible), I (Intermediate), and R
577 (Resistant) categories was performed according to the 2020 edition of EUCAST
578 recommendation on interpretation of the diameter of the zones of inhibition of
579 *Enterobacteriaceae*[37], and to the 2013 CA-SFM (Comité de l'Antibiogramme de la Société
580 Française de Microbiologie) standards for *Enterobacteriaceae*[38] for antibiotics for which
581 critical diameters are no longer reported in the latest published guidelines. *E. coli* CIP 76.24 (=
582 ATCC 25922) was used as a reference strain.

15

583

**DNA extraction and sequencing**

584

585 Genomic DNA was extracted at the NCPHL from subcultures inoculated with single bacterial
586 colonies and grown in nutrient agar (Oxoid, USA) at 37ºC overnight according to the
587 manufacturer instructions (Wizard® Genomic DNA Purification kit, Promega, UK). Genomic
588 DNA samples (derived from 10 environmental and 250 clinical samples, which includes the
589 120 samples sent to IP) were sent to the Wellcome Sanger Institute (WSI; Hinxton, UK) and
590 sequenced on the WSI sequencing pipeline (Figure S3) using the Illumina HiSeq platform X10
591 as previously described[28].

592 Two MDR *V. cholerae* strains were selected among the 22 held at the IP for long-read
593 sequencing. The first strain, CNRVC190243 (= YE-NCPHL-19014-PI), a 7PET *V. cholerae*
594 O1 strain was sequenced by Single-Molecule Real-Time (SMRT) sequencing (Pacific
595 Bioscience). The genomic DNA was prepared at the IP as follows: strain CNRVC190243 was
596 cultured in Brain-Heart-Infusion (BHI) broth (Difco) overnight at 37 °C with shaking (200
597 rpm—Thermo Scientific MaxQ 6800). Then, 100 µL of the overnight culture was inoculated
598 into a 10 ml BHI broth and cultured 2 hours at 37°C with shaking. After centrifugation, the
599 bacterial cells were processed as described previously[39], except that MaXtract High Density
600 columns (Qiagen) were used (instead of phase lock tubes) and DNA was resuspended in
601 molecular biology grade water (instead of 10 mM Tris pH 8.0). Library preparation and the
602 sequencing were performed by the GATC platform (Eurofins Genomics Europe Sequencing
603 GmbH; Konstanz, Germany) using their standard genomic library protocol and PacBio RS
604 sequencer. The second strain, CNRVC190247 (= YE-NCPHL-18035-PI), a non-O1/non-O139
605 *V. cholerae* strain further characterized as ST 555, was sequenced using the MinION nanopore
606 sequencer (Oxford Nanopore Technologies). Genomic DNA was prepared at the IP as follows:
607 strain CNRVC19247 was cultured in alkaline nutrient agar (casein meat peptone E2 from
608 Organotechnie, 20 g; sodium chloride from Sigma, 5 g; Bacto agar from Difco, 15g; distilled
609 water to 1 L; adjusted to pH 8.4; autoclaved at 121°C for 15 min) overnight at 37 °C. A few
610 isolated colonies of the overnight culture were inoculated into a 20 ml of Brain-Heart-Infusion
611 (BHI) broth and cultured until a final $OD_{600}=0.8$ at 37°C with shaking (200 rpm). After
612 centrifugation, the bacterial cells were processed as described above. The library was prepared
613 according to the instructions of the "Native barcoding genomic DNA (with EXP-NBD104,
614 EXP-NBD114, and SQK-LSK109)" procedure provided by Oxford Nanopore Technology.
615 The sequencing was then performed using a R9.4.1 flow cell on the MinION Mk1C apparatus
616 (Oxford Nanopore Technologies). The genomes of 21/22 strains cultivated at the IP (all but
617 CNRVC190251, which was isolated later; Table S2) were also sequenced using Illumina short-
618 read technology at the IP using the equipment and services of the iGenSeq platform at the
619 Institut du cerveau et de la moëlle épinière (Paris, France) from genomic DNA extracted with
620 the Maxwell 16-cell DNA purification kit (Promega) in accordance with the manufacturer's
621 recommendations.

622

**Genome assembly and annotation**

623

624 The 260 sequencing read sets produced at the WSI (Figure S3) were processed with the WSI
625 Pathogen Informatics pipeline[40]: quality of sequencing runs was assessed based on quality of
626 mapping of 10% reads to the genome of reference strain N16961 (GenBank Assembly

16

627  accession GCA_900205735.1) using the Burrows-Wheeler Aligner (BWA)[41]; read sets passed
628  the check if at least 80% bases were mapped after clipping, the base and indel error rate were
629  smaller than 0.02, and less than 80% of the insert sizes fell within 25% of the most frequent
630  size. Contamination was assessed manually based on Kraken classification of reads using the
631  standard WSI Pathogen reference database, which contains all viral, archaeal and bacterial
632  genomes and the mouse and human reference published in the RefSeq database as of the 21st
633  May 2015 (Table S3). Sequences were assembled *de novo* into contigs as described
634  previously[42], using SPAdes v3.10.0 as the core assembler[43]. Poor assemblies were filtered out
635  if differing of more than 20% from the expected genome size of 4.2 Mb, or when more than
636  10% of reads were assigned by Kraken to another organism than *V. cholerae* (notably including
637  the *Vibrio* phage ICP1) or to synthetic constructs, or were unclassified. This led to the omission
638  of 28 genome assemblies, resulting in 232 high-quality assembled genomes. The genome of
639  strains CNRVC190243 and CNRVC190247 were assembled based on long and short reads
640  using a hybrid approach with UniCycler[44] v0.4.7 and v0.4.8, respectively, using pilon[45] v1.23
641  for the polishing step, to produce high-quality reference sequences comprised of both
642  chromosomes and, for strain CNRVC190243, of an additional plasmid, pCNRVC190243. New
643  genomes were annotated with Prokka version v1.5.0[46].

**Contextual genomic data (882 "assembled *V. cholerae* genomes" dataset)**

646  To provide phylogenetic context, we also included in this analysis previously published
647  genome sequences from a globally representative set of isolates. We first gathered genome
648  assemblies generated at the WSI using the pipeline described above based on previously
649  published short reads sets from *V. cholerae* isolates belonging to sublineage T13 of 7PET Wave
650  3 (7PET-T13) and from strains isolated in close spatio-temporal context i.e. within a decade in
651  Africa and South Asia (where the ancestor of T13 is thought to originate[7]). These include all
652  42 Yemen 2016-2017 isolates[7], 103 recent isolates from East Africa including from Kenya[7],
653  Tanzania[47], Uganda[48] and Zimbabwe[15] and 74 isolates from South Asia [49]. In addition, we
654  included genomes spanning the wider diversity of *V. cholerae*, including all 119 genomes from
655  China[18], as well as 312 genomes from the collections of contextual genomes used in previous
656  studies[7,28]. Together with the 232 Yemen 2018-2019 isolate genome assemblies (see above),
657  our final dataset consisted of 882 assembled *V. cholerae* genomes (Table S4; Figure S3).

**Identification and typing of mobile genetic elements, virulence factors, AMR genes and
anti-phage defense systems**

661  The presence of AMR genes, plasmid replicon regions or virulence factors were predicted using
662  Abricate[50], searching the reference databases NCBI AMR+[51], Plasmidfinder[52] or VFDB[53],
663  respectively. BLASTN[54] (v2.7.1+, with default parameters) was used to identify known mobile
664  genetic elements (MGEs): the SXT/ICE ICE*Vch*Ind5 (GenBank accession GQ463142.1);
665  ICP1-like vibriophages ICP1_VMJ710 and ICP1_2012_A (GenBank accessions MN402506.2
666  and MH310936.1, respectively)[55] and the ICP1-like vibriophage YE-NCPHL-19021, which
667  genome was the only assembled contig from the reads obtained from sample YE-NCPHL-
668  19021 (this study; Genbank accession MW911613.1); the IncC-type plasmid
669  pCNRVC190243, obtained from the hybrid assembly of strain CNRVC190243 described
670  above (this study; ENA sequence accession OW443149.1); the MDR pseudo-compound

17

671 transposon (PCT) Yem*Vch*MDRI, extracted from this plasmid (positions 16,442 to 36,862);
672 PICI-like elements (PLE) 1, 2 and 3 (GenBank accessions KC152960.1, KC152961.1,
673 MF176135.1)[56,57]. Absence of elements was verified at the read level as described below.
674 Sequences similar to the reference sequences of the plasmid pCNRVC190243, the MDR PCT
675 Yem*Vch*MDRI and the ICP1-like phage genome YE-NCPHL-19021 were also searched in a
676 database of 661,405 genome assemblies [58] using a *k*-mer-based COBS index [59]; alignment of
677 best matches were further characterized using BLASTN. We typed the conjugation apparatus
678 of pCNRVC190243 with CONJScan [60] on the Pasteur Institute Galaxy server (Galaxy Version
679 1.0.5+galaxy0). We searched for presence of CRISPR-Cas arrays using MacSyFinder[61] v1.0.5
680 with default parameters and the built-in Cas system reference database; genomes positive for
681 Cas systems were further analysed with CRISPRCasFinder [62] on the Pasteur Institute Galaxy
682 server to retrieve CRISPR arrays.
683
684 **Prediction of serotype, serogroup and multi-locus sequene type**
685 To predict the antigenic serogroup, we screened the assemblies against a custom reference
686 database using Hamburger[63]. In brief, a database was constructed by selecting flanking and
687 marker genes for the operon encoding the *V. cholerae* O-antigen, with representative genes for
688 both O1 and non-O1 serogroups included (Table S15). Gene sequences were individually
689 aligned using Clustal Omega (version 1.2.4), prior to HMM construction with HMMER
690 (version 3.2.1) and concatenation of the HMM alignments. Assemblies were screened against
691 this database using Hamburger (version 836a77c)[64] to identify the operon, and genetic structure
692 was compared across the assemblies and references to designate serogroups. The HMMER
693 database is available online at https://figshare.com/s/5dd21a52f0d5a39a670f (doi:
694 10.6084/m9.figshare.19575148).
695 For O1 serotype prediction (Inaba or Ogawa), we used a combination of approaches including
696 BLASTN search against the 882 assembled *V. cholerae* genomes (as described above) and
697 ARIBA (v2.14.6+, with default parameters)[65] to screen the sequencing read sets against the
698 *wbeT* gene sequence from strain NCTC 9420 (positions 311,049-311,909 of GenBank
699 accession CP013319.1) as a reference, as previously described[28]. Multi-locus sequence typing
700 (MLST) of non-7PET isolates was conducted on PubMLST.org[66] under the non-O1/non-O139
701 *V. cholerae* seven-gene typing scheme.
702
703 **Identification of single nucleotide variants (456 "mapped 7PET genomes" and 33**
704 **"mapped *Vc*D genomes" datasets)**
705 For variant calling, Illumina short reads were mapped against the novel reference genomes
706 from strains CNRVC190243 and CNRVC190247, or the in-house MGE database described
707 above. We mapped all 260 short read sets from 2018-2019 Yemeni isolates sequenced at the
708 WSI, including those 28 read sets which assembly showed low coverage or appeared
709 contaminated with phage genomes (Table S3), so to recover variation data evidenced at the
710 read level, provided reads were mapped at a sufficient depth (see below). We also mapped read
711 sets from the 21 strains sequenced at the IP, and from contextual isolates of the 7PET-T13
712 sublineage and close relatives (see "Contextual genomic data"), for a total of 468 mapped
713 genomes. Reads were trimmed with Trimmomatic, mapped to both CNRVC190243 reference
714 chromosomes with BWA-MEM and the IncC plasmid pCNRVC190243. Mapped genomes

715 with an average read depth below 5x over the two chromosomes were deemed of insufficient
716 read depth and were excluded (12 read sets mapped to CNRVC190243, all from this study and
717 generated at WSI, were excluded for a final set of 456 mapped 7PET genomes [Table S5]; no
718 read set mapped to CNRVC190247 was excluded). We used the software suite
719 samtools/bcftools[67] v1.9 to call single nucleotide variants with a minimum coverage of 10x
720 read depth; see custom script 'map_yemen_reads2MGEs.sh'[68] for a detailed description of the
721 parameters used. Resulting consensus sequences were combined into a whole-genome
722 alignment, which was processed with snp-sites [69] to produce a single nucleotide polymorphism
723 (SNP) alignment.
724 Overall genome similarity was assessed by computing SNP distances based on the above
725 alignments using the function 'dist.dna' from the R package 'ape'[70], and average nucleotide
726 identity (ANI, accounting for unaligned regions) was computed using fastANI[71] v1.3 with
727 default parameters.
728
729 **Phylogenetic inference**
730 The Pantagruel pipeline[72] was used to infer a maximum-likelihood (ML) "core-genome tree"
731 using the "-S" option and otherwise default parameters. Briefly, 291 single-copy core-genome
732 genes (with expected high degree of sequence conservation and relatively low prevalence of
733 HGT compared to other core genes) were extracted from the 882 assembled *V. cholerae*
734 genomes, their alignments were concatenated and the resulting supermatrix was reduced to its
735 37,170 polymorphic positions, from which a ML tree was computed from RAxML v8.2.11[73]
736 (model ASC_GTRGAMMAX using Stamatakis' ascertainment bias correction; one starting
737 parsimony tree; 200 rapid bootstraps for estimating branch supports); supporting
738 supplementary data are available on Figshare at
739 https://figshare.com/s/3fe31c131b00a2a08bb9. Phylogenies were also inferred from whole-
740 genome alignments of the concatenated consensus sequences of both chromosomes from the
741 SNP alignment of the 456 mapped 7PET genomes and 33 mapped *Vc*D genomes. These
742 alignments contained 2,092 and 91,312 polymorphic positions, respectively, and were used as
743 input to RAxML-NG v1.0.1[74] to build the ML "mapped genome trees" using the following
744 options: "all --tree pars{10} --bs-trees 200 --model GTR+G4+ASC_STAM". Alternative
745 topologies were compared using RAxML-NG option "--sitelh" to generate per-site likelihood
746 values and the 'SH.test' function from the 'phangorn' R package[75] to test hypotheses.
747
748 The 882 assembled *V. cholerae* core-genome tree was rooted using the clade of sequences
749 identified as *V. paracholerae*[14] as an outgroup. The remaining part of the tree (*V. cholerae*
750 *sensu stricto*) was subdivided into clades named *Vc*A to *Vc*K based on visual examination with
751 the aim to coincide with previously described lineages such as 7PET, Gulf Coast, etc. or based
752 on balanced subdivisions of the tree diversity. *Vc*H, corresponding to the 7PET lineage, was
753 further subdivided into clades of even depth, named Subclades H.1 to H.9. The 456 mapped
754 7PET genomes were similarly classified into clusters based on the tree topology, with genomes
755 assigned to subclades named *Vc*H.5, *Vc*H.6, *Vc*H.8 or *Vc*H.9 (according to their position in the
756 882 assembled *V. cholerae* core-genome tree). Genomes belonging to *Vc*H.9, which
757 corresponds to the 7PET-T13 sublineage, were further separated into *Vc*H.9.a to *Vc*H.9.h,
758 based on visual examination of the tree structure and aiming to maximise uniformity of the

19

759    spatio-temporal metadata associated to genomes in each cluster; clusters correspond to clades,
760    either entirely or at the exclusion of another cluster included in the clade i.e. genome clusters
761    can emerge from each other. Final trees for the mapped genome datasets were rooted manually
762    according to the branching pattern in the 882 assembled *V. cholerae* core-genome tree, which
763    diversity encompases that of the mapped genome trees.

764

765    From a subset of the 456 mapped 7PET genome alignments (n=335) corresponding to *Vc*H.9,
766    a recombination-free phylogeny was inferred using ClonalFrameML v1.11[76] with default
767    parameters and using the ML mapped genome tree (restricted to the *Vc*H.9 genome tips) as a
768    starting tree. BactDating[77] v1.1 was then used to estimate a timed phylogeny (using 100,000
769    Monte-Carlo Markov chain iterations and otherwise default parameters) of the Yemen 2016-
770    2019 genomes and relatives using the ClonalFrameML tree and day-resolved dates as input;
771    median day of the year of isolation was used for isolates where these data were missing. Three
772    independent chains were run from different random seeds and yielded close results.

773

774    Supporting data for phylogenetic analyses of the 882 assembled *V. cholerae*, 456 mapped
775    7PET genomes and 33 mapped *Vc*D genomes are avaible on Figshare repository at
776    https://figshare.com/s/3fe31c131b00a2a08bb9 (doi: 10.6084/m9.figshare.16611823),
777    https://figshare.com/s/4d83a32cce78a52b413e (doi: 10.6084/m9.figshare.16595999) and
778    https://figshare.com/s/0be28064870c811120c5 (doi: 10.6084/m9.figshare.18304961),
779    respectively.

780

781    **Correlation of spatio-temporal and phylogenetic distances**
782    GPS data associated to the site of sample collection (health centres) were used to compute
783    spatial geodetic distances using R script 'gps_coords.r'[78,79]. Temporal distances were computed
784    from the difference between day of collection (only available for 2018 and 2019 Yemen
785    isolates). Phylogenetic distances were computed from the mapped genome tree using the
786    function 'cophenetic' from the core R package 'stats'[80]. Spatial, temporal and phylogenetic
787    distances were compared using a Monte-Carlo approximation of the Mantel test as implement
788    in the 'mantel.randtest' function from the R package 'ade4'[81], using 100,000 permutations to
789    compute the simulated *p*-value. Maps showing the distribution of genomes clusters over the
790    Yemen territory and in the region of Sana'a were obtained using QGIS 3.16.3 and the
791    QuickOSM API to retrieve OpenStreetMap data, specifically level 4 administrative boundaries
792    (governorates) in Yemen (last accessed 11 February 2021).

793

794    **Clade-specific SNPs and pangenome analysis**
795    The synteny-aware pangenome pipeline Panaroo[82] (v1.2.3) was run on the 882 assembled *V.*
796    *cholerae* genome set with the option "--clean-mode strict" and default parameters otherwise.
797    In parallel, a combined VCF file containing information on all SNP variation within the 456
798    mapped genome set was obtained using the 'bcftools merge' command. To identify clade-
799    specific SNPs and accessory gene presence/absence patterns, we used custom R scripts[68] to
800    compare the combined VCF file and the gene presence/absence table output of Panaroo,
801    respectively, to the mapped genome tree. Based on lists of genomes assigned to various clades
802    and clusters (see Results), we identify SNPs or accessory genes that are specific of a focus

20

803     clade in contrast to a background group or a sister clade, considering the contrast significant
804     when the Bonferonni-corrected p-value is below 0.05 and when the frequency of an allele is
805     above 0.8 in the focus clade and below 0.2 in the background clade, or conversely. Pangenome
806     analysis files are available at https://figshare.com/s/675d2a9e424ad4f11646 (doi:
807     10.6084/m9.figshare.19519105). Putative anti-phage defense systems were searched by testing
808     correlation of presence/absence patterns between ICP1-like phage and each pangenome gene
809     cluster; only associations with Pearson correlation coefficients lower than -0.9 or greater than
810     0.9 and p-values lower than $10^{-5}$ were retained as significant.
811

812 **Data availability**
813     Novel genomic data are available from the ENA/NCBI/DDBJ short read archive under the
814     BioProject PRJEB34436. Four of the resulting assemblies comprised a single 123-kb contig
815     corresponding to the ICP1-like phage; these assemblies were deemed uncontaminated and
816     complete ICP1-like phage genomes and were deposited to GenBank under the accessions
817     MW911612-MW911615. Complete hybrid genome assemblies for reference strains
818     CNRVC019243 and CNRVC019247 were deposited to the ENA under the BioProject
819     acessions PRJEB52123 and PRJEB47951 (Assemblies GCA_937000105 and
820     GCA_937000115), respectively. Suplementary data are available online on the Figshare
821     repository, under the following digital object ientifiers (doi):
822     https://doi.org/10.6084/m9.figshare.16595999, https://doi.org/10.6084/m9.figshare.16611823,
823     https://doi.org/10.6084/m9.figshare.18304961, https://doi.org/10.6084/m9.figshare.19097111,
824     https://doi.org/10.6084/m9.figshare.19519105, https://doi.org/10.6084/m9.figshare.19575148.
825
826

840

841 **References**
842

843     1.   UNHCR Yemen 2021 Country Operational Plan. *UNHCR Operational Data Portal*

844        *(ODP)* https://data2.unhcr.org/en/documents/details/85850 (2021).

845   2.  Dureab, F. *et al.* Assessment of electronic disease early warning system for improved

846       disease surveillance and outbreak response in Yemen. *BMC Public Health* **20**, 1422

847       (2020).

848   3.  WHO-EMRO. Cholera outbreaks. *World Health Organization - Regional Office for the*

849       *Eastern Mediterranean* http://www.emro.who.int/health-topics/cholera-outbreak/cholera-

850       outbreaks.html (2021).

851   4.  WHO. Health workers in Yemen reach more than 306,000 people with cholera vaccines

852       during four-day pause in fighting – WHO, UNICEF. *World Health Organization*

853       https://www.who.int/news/item/05-10-2018-health-workers-in-yemen-reach-more-than-

854       306-000-people-with-cholera-vaccines-during-four-day-pause-in-fighting-who-unicef

855       (2018).

856   5.  Federspiel, F. & Ali, M. The cholera outbreak in Yemen: lessons learned and way

857       forward. *BMC Public Health* **18**, 1338 (2018).

858   6.  Mutreja, A. *et al.* Evidence for several waves of global transmission in the seventh

859       cholera pandemic. *Nature* **477**, 462–465 (2011).

860   7.  Weill, F.-X. *et al.* Genomic insights into the 2016–2017 cholera epidemic in Yemen.

861       *Nature* **565**, 230 (2019).

862   8.  Bai, Z. G. *et al.* Azithromycin versus penicillin G benzathine for early syphilis. *Cochrane*

863       *Database Syst Rev* CD007270 (2012) doi:10.1002/14651858.CD007270.pub2.

864   9.  Rabaan, A. A. Cholera: an overview with reference to the Yemen epidemic. *Front Med*

865       **13**, 213–228 (2019).

866   10. Global Task Force on Cholera Control. Technical note on the use of antibiotics for the

867       treatment and control of cholera. https://www.gtfcc.org/wp-

868       content/uploads/2019/10/gtfcc-technical-note-on-use-of-antibiotics-for-the-treatment-of-

869       cholera.pdf (2018).

870  11. Camacho, A. *et al.* Cholera epidemic in Yemen, 2016–18: an analysis of surveillance

871      data. *The Lancet Global Health* **6**, e680–e690 (2018).

872  12. Ambrose, S. J., Harmer, C. J. & Hall, R. M. Compatibility and entry exclusion of IncA

873      and IncC plasmids revisited: IncA and IncC plasmids are compatible. *Plasmid* **96–97**, 7–

874      12 (2018).

875  13. De, R. Mobile Genetic Elements of Vibrio cholerae and the Evolution of Its

876      Antimicrobial Resistance. *Frontiers in Tropical Diseases* **2**, 7 (2021).

877  14. Islam, M. T. *et al.* Population analysis of Vibrio cholerae in aquatic reservoirs reveals a

878      novel sister species (Vibrio paracholerae sp. nov.) with a history of association with

879      human infections. *bioRxiv* 2021.05.05.442690 (2021) doi:10.1101/2021.05.05.442690.

880  15. Mashe, T. *et al.* Highly Resistant Cholera Outbreak Strain in Zimbabwe. *New England*

881      *Journal of Medicine* **383**, 687–689 (2020).

882  16. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new

883      developments. *Nucleic Acids Research* **47**, W256–W259 (2019).

884  17. Greig, D. R. *et al.* Evaluation of Whole-Genome Sequencing for Identification and

885      Typing of Vibrio cholerae. *Journal of Clinical Microbiology* **56**, (2018).

886  18. Wang, H. *et al.* Genomic epidemiology of Vibrio cholerae reveals the regional and global

887      spread of two epidemic non-toxigenic lineages. *PLOS Neglected Tropical Diseases* **14**,

888      e0008046 (2020).

889  19. Spagnoletti, M. *et al.* Acquisition and Evolution of SXT-R391 Integrative Conjugative

890      Elements in the Seventh-Pandemic Vibrio cholerae Lineage. *mBio* **5**, e01356-14 (2014).

891  20. Harmer, C. J., Pong, C. H. & Hall, R. M. Structures bounded by directly-oriented

892      members of the IS26 family are pseudo-compound transposons. *Plasmid* **111**, 102530

893      (2020).

23

21. Bonnin, R. A. *et al.* PER-7, an Extended-Spectrum β-Lactamase with Increased Activity toward Broad-Spectrum Cephalosporins in Acinetobacter baumannii. *Antimicrobial Agents and Chemotherapy* **55**, 2424–2427 (2011).

22. Toleman, M. A., Bennett, P. M. & Walsh, T. R. ISCR Elements: Novel Gene-Capturing Systems of the 21st Century? *Microbiology and Molecular Biology Reviews* **70**, 296–316 (2006).

23. Yu, L. *et al.* Multiple Antibiotic Resistance of Vibrio cholerae Serogroup O139 in China from 1993 to 2009. *PLOS ONE* **7**, e38633 (2012).

24. Wang, R. *et al.* IncA/C plasmids harboured in serious multidrug-resistant Vibrio cholerae serogroup O139 strains in China. *International Journal of Antimicrobial Agents* **45**, 249–254 (2015).

25. Opazo, A. *et al.* Plasmid-encoded PER-7 β-lactamase responsible for ceftazidime resistance in Acinetobacter baumannii isolated in the United Arab Emirates. *Journal of Antimicrobial Chemotherapy* **67**, 1619–1622 (2012).

26. Jaskólska, M., Adams, D. W. & Blokesch, M. Two defence systems eliminate plasmids from seventh pandemic Vibrio cholerae. *Nature* 1–7 (2022) doi:10.1038/s41586-022-04546-y.

27. Ceccarelli, D., Salvia, A. M., Sami, J., Cappuccinelli, P. & Colombo, M. M. New Cluster of Plasmid-Located Class 1 Integrons in Vibrio cholerae O1 and a dfrA15 Cassette-Containing Integron in Vibrio parahaemolyticus Isolated in Angola. *Antimicrobial Agents and Chemotherapy* (2006) doi:10.1128/AAC.01310-05.

28. Dorman, M. J. *et al.* Genomics of the Argentinian cholera epidemic elucidate the contrasting dynamics of epidemic and endemic Vibrio cholerae. *Nature Communications* **11**, 4918 (2020).

918    29. Haley, B. J. *et al.* Genomic and Phenotypic Characterization of Vibrio cholerae Non-O1

919        Isolates from a US Gulf Coast Cholera Outbreak. *PLOS ONE* **9**, e86264 (2014).

920    30. Weill, F.-X. *et al.* Genomic history of the seventh pandemic of cholera in Africa. *Science*

921        **358**, 785–789 (2017).

922    31. Dhabaan, G., Chahin, A., Buhaish, A. & Shorman, M. COVID-19 pandemic in Yemen: A

923        questionnaire based survey, what do we know? *The Journal of Infection in Developing*

924        *Countries* **14**, 1374–1379 (2020).

925    32. Madico, G. *et al.* Active surveillance for Vibrio cholerae O1 and vibriophages in sewage

926        water as a potential tool to predict cholera outbreaks. *Journal of Clinical Microbiology*

927        **34**, 2968–2972 (1996).

928    33. CDC. Laboratory Methods for the Diagnosis of Vibrio cholerae. *Centers for Disease*

929        *Control and Prevention* https://www.cdc.gov/cholera/pdf/laboratory-methods-for-the-

930        diagnosis-of-vibrio-cholerae-chapter-4.pdf (2021).

931    34. Weinstein, M. P. *M100: Performance Standards for Antimicrobial Susceptibility Testing,*

932        *31st Edition*. (Clinical & Laboratory Standards Institute).

933    35. Dodin, A. & Fournier, J. M. Laboratory methods for the diagnosis of cholera vibrio and

934        other vibrios. in *Diagnosis of the cholera vibrio.* 59–82 (Institut Pasteur, 1992).

935    36. CASFM & EUCAST. Recommandations 2020, V.1.1 (Avril). (2020).

936    37. EUCAST. The European Committee on Antimicrobial Susceptibility Testing. Breakpoint

937        tables for interpretation of MICs and zone diameters. Version 10.0. (2020).

938    38. CASFM. Recommandations 2013. (2013).

939    39. von Mentzer, A. *et al.* Long-read-sequenced reference genomes of the seven major

940        lineages of enterotoxigenic Escherichia coli (ETEC) circulating in modern time. *Sci Rep*

941        **11**, 9256 (2021).

25

942    40. WSI Pathogen Informatics. *vr-codebase Github repository*. (Pathogen Informatics,

943        Wellcome Sanger Institute, 2022).

944    41. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.

945        *arXiv* (2013) doi:1303.3997v1 [q-bio.GN].

946    42. Page, A. J. *et al.* Robust high-throughput prokaryote de novo assembly and improvement

947        pipeline for Illumina data. *Microbial Genomics* **2**, (2016).

948    43. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its Applications

949        to Single-Cell Sequencing. *Journal of Computational Biology* **19**, 455–477 (2012).

950    44. Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial

951        genome assemblies from short and long sequencing reads. *PLOS Computational Biology*

952        **13**, e1005595 (2017).

953    45. Walker, B. J. *et al.* Pilon: An Integrated Tool for Comprehensive Microbial Variant

954        Detection and Genome Assembly Improvement. *PLOS ONE* **9**, e112963 (2014).

955    46. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–

956        2069 (2014).

957    47. Kachwamba, Y. *et al.* Genetic Characterization of Vibrio cholerae O1 isolates from

958        outbreaks between 2011 and 2015 in Tanzania. *BMC Infectious Diseases* **17**, 157 (2017).

959    48. Bwire, G. *et al.* Molecular characterization of Vibrio cholerae responsible for cholera

960        epidemics in Uganda by PCR, MLVA and WGS. *PLOS Neglected Tropical Diseases* **12**,

961        e0006492 (2018).

962    49. Morita, D. *et al.* Whole-Genome Analysis of Clinical Vibrio cholerae O1 in Kolkata,

963        India, and Dhaka, Bangladesh, Reveals Two Lineages of Circulating Strains, Indicating

964        Variation in Genomic Attributes. *mBio* **11**, (2020).

965    50. Seemann, T. *ABRicate*. (2021).

966   51. Feldgarden, M. *et al.* Validating the AMRFinder Tool and Resistance Gene Database by

967       Using Antimicrobial Resistance Genotype-Phenotype Correlations in a Collection of

968       Isolates. *Antimicrobial Agents and Chemotherapy* **63**, e00483-19.

969   52. Carattoli, A. *et al.* In Silico Detection and Typing of Plasmids using PlasmidFinder and

970       Plasmid Multilocus Sequence Typing. *Antimicrobial Agents and Chemotherapy* **58**,

971       3895–3903 (2014).

972   53. Chen, L., Zheng, D., Liu, B., Yang, J. & Jin, Q. VFDB 2016: hierarchical and refined

973       dataset for big data analysis—10 years on. *Nucleic Acids Res* **44**, D694–D697 (2016).

974   54. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421

975       (2009).

976   55. Angermeyer, A., Das, M. M., Singh, D. V. & Seed, K. D. Analysis of 19 Highly

977       Conserved Vibrio cholerae Bacteriophages Isolated from Environmental and Patient

978       Sources Over a Twelve-Year Period. *Viruses* **10**, 299 (2018).

979   56. Seed, K. D., Lazinski, D. W., Calderwood, S. B. & Camilli, A. A bacteriophage encodes

980       its own CRISPR/Cas adaptive response to evade host innate immunity. *Nature* **494**, 489–

981       491 (2013).

982   57. O'Hara, B. J., Barth, Z. K., McKitterick, A. C. & Seed, K. D. A highly specific phage

983       defense system is a conserved feature of the Vibrio cholerae mobilome. *PLoS Genet* **13**,

984       e1006838 (2017).

985   58. Blackwell, G. A. *et al.* Exploring bacterial diversity via a curated and searchable snapshot

986       of archived DNA sequences. *PLOS Biology* **19**, e3001421 (2021).

987   59. Bingmann, T., Bradley, P., Gauger, F. & Iqbal, Z. COBS: A Compact Bit-Sliced

988       Signature Index. in *String Processing and Information Retrieval* (Springer International

989       Publishing, 2019).

27

990    60. Cury, J., Abby, S. S., Doppelt-Azeroual, O., Néron, B. & Rocha, E. P. C. Identifying

991        Conjugative Plasmids and Integrative Conjugative Elements with CONJscan. in

992        *Horizontal Gene Transfer: Methods and Protocols* (ed. de la Cruz, F.) 265–283 (Springer

993        US, 2020). doi:10.1007/978-1-4939-9877-7_19.

994    61. Abby, S. S., Néron, B., Ménager, H., Touchon, M. & Rocha, E. P. C. MacSyFinder: A

995        Program to Mine Genomes for Molecular Systems with an Application to CRISPR-Cas

996        Systems. *PLoS ONE* **9**, e110726 (2014).

997    62. Pourcel, C. *et al.* CRISPRCasdb a successor of CRISPRdb containing CRISPR arrays and

998        cas genes from complete genome sequences, and tools to download and query lists of

999        repeats and spacers. *Nucleic Acids Research* **48**, D535–D544 (2020).

1000   63. Williams, D. J. *et al. The phylogenomic landscape of the genus Serratia.*

1001       2022.01.11.475790 https://www.biorxiv.org/content/10.1101/2022.01.11.475790v1

1002       (2022) doi:10.1101/2022.01.11.475790.

1003   64. Williams, D. *hamburger.* (2021).

1004   65. Hunt, M. *et al.* ARIBA: rapid antimicrobial resistance genotyping directly from

1005       sequencing reads. *Microbial Genomics* **3**, (2017).

1006   66. Jolley, K. A., Bray, J. E. & Maiden, M. C. J. Open-access bacterial population genomics:

1007       BIGSdb software, the PubMLST.org website and their applications. (2018)

1008       doi:10.12688/wellcomeopenres.14826.1.

1009   67. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience* **10**, giab008

1010       (2021).

1011   68. Lassalle, F. *flass/yemenpaper Github repository.* (Wellcome Sanger Institute, 2022).

1012   69. Page, A. J. *et al.* SNP-sites: rapid efficient extraction of SNPs from multi-FASTA

1013       alignments. *Microbial Genomics,* **2**, e000056 (2016).

1014   70. Paradis, E. *Analysis of Phylogenetics and Evolution with R.* (Springer New York, 2012).

1015    71. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High

1016    throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.

1017    *Nature Communications* **9**, 5114 (2018).

1018    72. Lassalle, F., Jauneikaite, E., Veber, P. & Didelot, X. Automated reconstruction of all

1019    gene histories in large bacterial pangenome datasets and search for co-evolved gene

1020    modules with Pantagruel. *bioRxiv* 586495 (2019) doi:10.1101/586495.

1021    73. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of

1022    large phylogenies. *Bioinformatics* btu033 (2014) doi:10.1093/bioinformatics/btu033.

1023    74. Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAxML-NG: a fast,

1024    scalable and user-friendly tool for maximum likelihood phylogenetic inference.

1025    *Bioinformatics* doi:10.1093/bioinformatics/btz305.

1026    75. Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592–593 (2011).

1027    76. Didelot, X. & Wilson, D. J. ClonalFrameML: Efficient Inference of Recombination in

1028    Whole Bacterial Genomes. *PLoS Comput Biol* **11**, e1004041 (2015).

1029    77. Didelot, X., Croucher, N. J., Bentley, S. D., Harris, S. R. & Wilson, D. J. Bayesian

1030    inference of ancestral dates on bacterial phylogenetic trees. *Nucleic Acids Research* **46**,

1031    e134–e134 (2018).

1032    78. Lassalle, F. *flass/microbiomes Github repository*. (2018).

1033    79. Lassalle, F. *et al.* Oral microbiomes from hunter-gatherers and traditional farmers reveal

1034    shifts in commensal balance and pathogen load linked to diet. *Molecular Ecology* **27**,

1035    182–195 (2018).

1036    80. R Core Team. *R: A Language and Environment for Statistical Computing.* (R Foundation

1037    for Statistical Computing, 2020).

1038    81. Dufour, A.-B. & Dray, S. The ade4 Package: Implementing the Duality Diagram for

1039    Ecologists. *Journal of Statistical Software* **22**, (2007).

1040    82. Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo

1041        pipeline. *Genome Biology* **21**, 180 (2020).

1042    83. Sullivan, M. J., Petty, N. K. & Beatson, S. A. Easyfig: a genome comparison visualizer.

1043        *Bioinformatics* **27**, 1009–1010 (2011).

1044    84. Hoshino, K. *et al.* Development and evaluation of a multiplex PCR assay for rapid

1045        detection of toxigenic Vibrio cholerae O1 and O139. *FEMS Immunology & Medical*

1046        *Microbiology* **20**, 201–207 (1998).

1047    85. Nonaka, L. & Suzuki, S. New Mg2+-Dependent Oxytetracycline Resistance Determinant

1048        Tet 34 in Vibrio Isolates from Marine Fish Intestinal Contents. *Antimicrobial Agents and*

1049        *Chemotherapy* **46**, 1550–1552 (2002).

1050    86. Roberts, M. C. Update on acquired tetracycline resistance genes. *FEMS Microbiology*

1051        *Letters* **245**, 195–203 (2005).

1052    87. Seed, K. D. *et al.* Evidence of a Dominant Lineage of Vibrio cholerae-Specific Lytic

1053        Bacteriophages Shed by Cholera Patients over a 10-Year Period in Dhaka, Bangladesh.

1054        *mBio* **2**, (2011).

1055    88. Nelson, E. J. *et al.* Gold Standard Cholera Diagnostics Are Tarnished by Lytic

1056        Bacteriophage and Antibiotics. *Journal of Clinical Microbiology* **58**, (2020).

1057    89. LeGault, K. N. *et al.* Temporal shifts in antibiotic resistance elements govern phage-

1058        pathogen conflicts. *Science* **373**, (2021).

1059
1060
1061
1062

**Table 1. Number of *V. cholerae* isolate genomes from Yemen by year and phylogenetic lineage**

| Year | Total | Clades *Vpc*[1] | *Vc*D[2] | *Vc*K[1] | Clusters *Vc*H/H.9[3] | H.9.e | H.9.f | H.9.g | H.9.h | n.d.[4] |
|---|---|---|---|---|---|---|---|---|---|---|
| 2016 | 8 | | | | 8 | 7 | 1 | | | |
| 2017 | 34 | | | | 34 | 29 | 5 | | | |
| 2018 | 112 | | 17 | | 87 | 3 | | 6 | 78 | 8 |
| 2019 | 169 | 1 | 4 | 2 | 151 | | | 150 | 1 | 11 |
| Total | 323 | 1 | 21 | 2 | 280 | 39 | 6 | 156 | 79 | 19 |

[1] Assigned based on the "882 assembled *V. cholerae* genomes" dataset

[2] Assigned based on the "33 mapped *Vc*D genomes" dataset

[3] Assigned based on the 456 "mapped 7PET genomes"

[4] Poor quality genome data or no coverage of the bacterial genomes (e.g. in case of complete contamination by ICP1 virus genome).

**Figure Legends**


**Figure 1: Phylogenetic diversity of *Vibrio cholerae* isolates from Yemen**
Maximum-likelihood phylogeny of 882 assembled *V. cholerae* genomes based on the 37,170 SNP sites from the concatenated alignments of 291 core genes. Low-diversity clades (*Vc*H and part of *Vc*K) are collapsed and marked by black stars. Clades are highlighted with background colours (legend key 1). Coloured rings outside the tree depict the match with previously described lineages (ring 2), the geographical origin of isolates at the level of continents (ring 3), and their year of isolation when from Yemen (ring 4). Presence of parts of the plasmid pCNRVC190243 are indicated by coloured circles (ring 5 in A): IncC plasmid backbone (light brown) and the MDR pseudo-compound transposon Yem*Vch*MDRI (dark brown); full circles indicate over 70% coverage in assemblies of the reference length, hollow circles indicate 30-70% coverage in assemblies and confirmed presence based on mapped reads, with even coverage over the MGE reference sequence, while half-circles represent heterogeneous presence in a collapsed clade. Tree plots were generated with iTOL v4[16] and adapted with Inkscape. The scale bar represents the number of nucleotide substitutions per variable site.


**Figure 2: Phylogenetic diversity and spatiotemporal distribution of *Vibrio cholerae* 7PET-T13 isolates (*Vc*H.9) from Yemen**
**A.** Subtree of the maximum-likelihood phylogeny of 456 7PET genomes mapped to reference *Vc*H.9 strain CNRVC190243 genome, including 335/456 genomes covering *Vc*H.9 (as defined in Figure S5), which corresponds to the 7PET-T13 sublineage and close South Asian relatives. The full tree containing the 456 genomes is available as supplementary material on Figshare (https://figshare.com/s/4d83a32cce78a52b413e; doi: 10.6084/m9.figshare.16595999) and was obtained based on 2,092 SNP sites from concatenated whole-chromosome alignments. Brown branches indicate the clade grouping all Yemeni 7PET-T13 isolates. Bootstrap support over 70% is indicated by white circles. Phylogenetic clusters within *Vc*H.9 are highlighted with background colours (legend key 1). Coded tracks outside the tree depict the serotype of isolates (ring 2) as predicted from genomic data, year of isolation when isolated in 2012 or later (ring 3), the governorate of isolation if in Yemen (ring 4). The presence of mobile genetic elements (MGEs) is indicated by coloured circles in the outermost track (ring 5): ICP1-like phage (pink), SXT/ICE ICE*Vch*Ind5 (blue), ICE*Vch*Ind5$^\Delta$ i.e. featuring the characteristic 10-kb deletion in the variable region III (green), IncC plasmid backbone (light brown) and the MDR pseudo-compound transposon Yem*Vch*MDRI (dark brown); filled and unfilled circles indicate different level of coverage in assemblies (see Figure 1 legend). The position of the reference sequence to which all other genomes were mapped to generate the alignment is labelled. The scale bar represents the number of nucleotide substitutions per site. **B.** Frequency of each phylogenetic subcluster among Yemen isolates per month since the onset of the Yemen outbreak. Where relevant, the cluster group is subdivided by the presence or absence of the IncC plasmid as indicated by the filled brown (present) or open (absence) circle on the right of the chart. The contribution of each governorate of isolation is indicated by the coloured portion of each bar. **C and D.** A map of Yemen governorates (C) and a focus on the Sana'a and Amanat Al Asimah governorates (inner and outer capital city; D), with dots corresponding to isolates, coloured by phylogenetic subcluster.

1125

1126 **Figure 3: Genetic organisation of the MDR pseudo-compound transposon Yem*Vch*MDRI**

1127 Antimicrobial resistance (AMR) genes are filled in black and labelled in boldface; genes

1128 encoding endonucleases transposases and other genes involved in genetic mobility are filled in

1129 grey. Genomic position is indicated by tickmarks every kilobase, in reference to the

1130 pCNRVC190243 plasmid coordinates.

**A**

Tree scale: 0.00001

5 4 3 2 1

Ref ▷

H.9.e
H.9.f
H.9.h
H.9.g

**1. Clusters within clade *Vc*H.9 / 7PET**
- *Vc*H.9.a
- *Vc*H.9.b
- *Vc*H.9.c
- *Vc*H.9.d
- *Vc*H.9.e
- *Vc*H.9.f
- *Vc*H.9.g
- *Vc*H.9.h

**2. Serotype**
- intact *wbeT* gene (predicted Ogawa serotype)

**3. Year of isolation**
- < 2012
- 2012
- 2013
- 2014
- 2015
- 2016
- 2017
- 2018
- 2019

**4. Yemen governorate**
- Al Bayda
- Al Mahwit
- Aljawf
- Amanat Al Asimah
- Amran
- Dhamar
- Hajjah
- Ibb
- Sana'a
- Taiz

**5. MGEs**
- ICP1-like phage
- ICEVchInd5
- ICEVchInd5Δ
- IncC plasmid backbone
- YemVchMDRI

**B**

*Vc*H.9.e

*Vc*H.9.f

*Vc*H.9.g

*Vc*H.9.h

Number of Isolates

Month, Year

Sep-2016, Nov-2016, Jan-2017, Mar-2017, May-2017, Jul-2017, Sep-2017, Nov-2017, Jan-2018, Mar-2018, May-2018, Jul-2018, Sep-2018, Nov-2018, Jan-2019, Mar-2019, May-2019, Jul-2019, Sep-2019, Nov-2019

**C**



**D**