1  **Genomic characteristics of recently recognized *Vibrio cholerae* El Tor lineages associated**
2  **with cholera in Bangladesh, 1991-2017**

3  **Authors**:

4  Md Mamun Monir[1], Talal Hossain[1], Masatomo Morita[2], Makoto Ohnishi[2], Fatema-Tuz Johura[1],
5  Marzia Sultana[1], Shirajum Monira[1], Tahmeed Ahmed[1], Nicholas Thomson[3], Haruo Watanabe[2],
6  Anwar Huq[4], Rita R. Colwell[4,5], Kimberley Seed[6], and Munirul Alam[1§].

7

8  **Authors affiliations**:

9  *1. icddr, b, Formerly International Centre for Diarrhoeal Disease Research, Bangladesh,*
10 *Dhaka, Bangladesh*
11
12 *2. National Institutes of Infectious Diseases (NIID), Tokyo, Japan*
13
14 *3. Sanger Institute, Cambridge, UK*
15
16
17 *4. Maryland Pathogen Research Institute, University of Maryland, USA*
18
19 *5. Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland, USA*
20
21 *6. University of California, Berkeley, USA*
22
23
24
25
26 [§]Corresponding author:  Munirul Alam, PhD

27 Mailing address:

28 Infectious Diseases Division (IDD)
29 International Centre for Diarrheal Disease Research, Bangladesh (icddr,b)
30 68, Shaheed Tajuddin Ahmed Sarani,
31 Mohakhali, Dhaka 1212,
32 Bangladesh
33 Tel:  +88-02-9840523-32 Ext. 2433/2490

34 Fax:  +88-02-8812529

35 E-Mail: munirul@icddrb.org

36  **Abstract: (Words 245)**

37  Comparative genomic analysis of *Vibrio cholerae* El Tor associated with endemic cholera in

38  Asia revealed two distinct lineages, one dominant in Bangladesh and the other in India.  An in

39  depth whole genome study of *V. cholerae* El Tor clinical strains isolated during endemic cholera

40  in Bangladesh (1991 – 2017) included reference genome sequence data obtained online. Core

41  genome phylogeny established using single nucleotide polymorphisms (SNPs) showed *V.*

42  *cholerae* El Tor strains comprised two lineages, BD-1 and BD-2, which, according to Bayesian

43  phylodynamic analysis, originated from paraphyletic group BD-0 around 1981. BD-1 and BD-2

44  lineages overlapped temporally but were negatively associated as causative agents of cholera

45  2004-2017. Genome wide association study (GWAS) revealed 140 SNPs and 31 indels, resulting

46  in gene alleles unique to BD-1 and BD-2. Regression analysis of root to tip distance and year of

47  isolation indicated early BD-0 strains at the base, whereas BD-1 and BD-2 subsequently emerged

48  and progressed by accumulating SNPs. Pangenome analysis provided evidence of gene

49  acquisition by both BD-1 and BD-2, of which six crucial proteins of known function were

50  predominant in BD-2. BD-1 and BD-2 diverged and have distinctively different genomic traits,

51  namely heterogeneity in VSP-2, VPI-1, mobile elements, toxin encoding elements, and total gene

52  abundance. In addition, the observed phage-inducible chromosomal island-like element (PLE1),

53  and SXT ICE elements (ICE$^{TET}$) in BD-2 presumably provided a fitness advantage for the

54  lineage to outcompete BD-1 as the etiological agent of the endemic cholera in Bangladesh, with

55  implications for global cholera epidemiology.

56  **Importance:** (150 words)

57  Cholera is a global disease with specific reference to the Bay of Bengal Ganges Delta where

58  *Vibrio cholerae* O1 El Tor, causative agent of the disease showed two circulating lineages, one

59  dominant in Bangladesh and the other in India. Results of in-depth genomic study of *V. cholerae*

60  associated with endemic cholera during the past 27 years (1991 – 2017) indicate emergence and

61  succession of the two lineages, BD-1 and BD-2, arising from a common ancestral paraphylatic

62  group, BD-0, comprising the early strains and short-term evolution of the bacterium in

63  Bangladesh. Among the two *V. cholerae* lineages, BD-2 supersedes BD-1 and is predominant in

64  the most recent endemic cholera in Bangladesh. The BD-2 lineage contained significantly more

65   SNPs and indels, and showed richness in gene abundance, including antimicrobial resistance

66   genes, gene cassettes, and PLE to fight against bacteriophage infection, acquired over time.

67   These findings have important epidemic implications at a global scale.

## Introduction

69   Cholera is a life threatening infectious diarrheal disease caused by *Vibrio cholerae* serogroups

70   O1 and O139 of the Gram-negative gammaproteobacteria (1, 2). The global incidence of cholera

71   is estimated to be 2.9 million cases annually with almost 95,000 deaths (3). In 2017, 34 countries

72   reported a total of 1,227,391 cases and 5,654 deaths (4). Seven cholera pandemics have been

73   recognized since 1817. However, limited information is available regarding the etiological agent

74   for the first five pandemics and no isolates of the causative agent are extant. The sixth pandemic,

75   and possibly those earlier were caused by *V. cholerae* O1 classical biotype, while the ongoing

76   seventh pandemic is caused by *V. cholerae* El Tor biotype and began with displacement of *V.*

77   *cholerae* classical biotype in Asia in 1961 (5). *V. cholerae* El Tor was isolated in Africa in the

78   1970s and Latin America in 1991 where for more than a century there had been no cholera

79   outbreaks (6). In 1992, a *V. cholerae* non-O1 strain designated *V. cholerae* O139 Bengal initiated

80   outbreaks of cholera in coastal areas of India and Bangladesh, and subsequently was isolated

81   from patients in several countries of Asia (2). *V. cholerae* El Tor continues to be the major

82   etiological agent of cholera worldwide.

83   The severe dehydrating diarrhea characteristic of cholera is associated with several factors,

84   including a toxin and several virulence genes involved in colonization and toxicity and their

85   coordinated expression (1). Cholera toxin (CT) is the virulence factor responsible for secretory

86   diarrhea of cholera and is encoded in the genome of a lysogenic CTX phage. *V. cholerae* El Tor

87   responsible for the current cholera pandemic harbors the CTX phage classical biotype variant,

88   and the ctxB$^{cla}$: ctxB genotype 1 (*ctxB*1) or *ctxB*7 (7). *V. cholerae* responsible for the current

89   cholera pandemic has become more virulent by undergoing several shifts in CTX genotype and

90   acquiring virulence-related gene islands (8). Integrative conjugative elements (ICEs) and

91   lysogenic phages are genetic elements that play an important role in the acquisition of virulence,

92   antimicrobial resistance, and heavy metal resistance, which important components of the

93   pathogenicity of *V. cholerae* (9, 10). Functions of these elements are important for the pathogen

94   to exert evolutionary advantage and variants can be used as markers of clonal expansion (1).

95　Acquisition of mobile genetic elements (MGEs) through horizontal gene transfer (HGT) and

96　propitious chromosomal mutations are significant landmarks for an evolving bacterium (11).

97　Whole-genome sequencing of *V. cholerae* El Tor strains associated with the seventh cholera

98　pandemic revealed three waves, suggesting independent but overlapping paths for the pathogen

99　to spread globally from the Bay of Bengal estuary where cholera has been endemic at least since

100　1961 but likely for centuries (5). Intercontinental transmission of *V. cholerae* has been proposed

101　for the 2010 outbreak in Haiti (12). Bangladesh borders on the Bay of Bengal and is considered

102　to be a hotspot of Asiatic cholera where -ca. 100,000 cases and 4,500 deaths are reported each

103　year (13). *V. cholerae* O1 responsible for endemic cholera in Bangladesh and India has been

104　found to have undergone genetic changes over time including acquisition of classical biotype

105　attributes in an El Tor background, thereby becoming more successful as a pathogen　(14, 15). A

106　recent whole-genome analysis of *Vibrio cholerae* El Tor strains isolated between 2009 and 2016

107　indicated two distinct lineages exist in Bengal (16). The objective of the study reported here was

108　to investigate *V. cholerae* endemic cholera strains isolated during 1991 to 2017 to understand

109　more completely about emergence and progression of the two lineages in Bangladesh. Virulence

110　and related genomic islands, including toxin and antimicrobial resistance genes differing

111　significantly among the *V. cholerae* El Tor lineages, were also investigated for potential

112　relevance to emergence of the lineages.

113　**Results**

114　**Phylogenetic analysis**

115　A total of 119 strains were included in the study and their genomes were sequenced using the

116　Illumina platform (MiSeq or HiSeq 2500 sequencer). In addition, 56 strains from our previous

117　study (16) and 17 genomes from the European Nucleotide Archive (17) were used, which are

118　representative of isolates from Bangladesh between 1991 and 2017 (see Table S1 Supplemental

119　Material). Paired-end reads of the 192 genomes were mapped to *V. cholerae* El Tor N16961

120　reference strain, a seventh-pandemic *V. cholerae* O1 El Tor (7PET) strain isolated in Bangladesh

121　in 1975 (18). A total of 1,298 single nucleotide polymorphisms (SNPs) and 413 indels (insertions

122　or deletions) were obtained and, after filtering indels, low call rate, and high-density SNPs, a

123　total of 893 high-quality SNPs were retained for further study. A phylogenetic analysis was

124  conducted to construct a tree based on the 893 high-quality SNPs to evaluate the genetic
125  diversity of the *Vibrio cholerae* O1 El Tor isolates from Bangladesh. A nested hierarchical
126  structure in the phylogenetic tree was observed, with all but four of the strains isolated between
127  1999 and 2017 clustering into two major clades, BD-1 (n=76) and BD-2 (n=105), shown in green
128  and red, respectively. The remaining strains formed paraphyletic group BD-0 (n=11) (Fig. 1A, in
129  blue). Except for three strains isolated in 2012 that formed a sub-clade, BD-0 consisted mostly of
130  strains isolated earlier between 1991 and 2000. Dates of isolation of common ancestors of the
131  lineages were inferred using Bayesian Markov chain Monte Carlo framework Bayesian
132  Evolutionary Analysis Sampling Trees (BEAST) (19) (see Fig. S1 supplementary material), and
133  a maximum clade credibility (MCC) tree was inferred from the posterior distribution of the best
134  fitting model using program TreeAnnotator tool of the BEAST software package. It was
135  estimated from the MCC tree that the most recent common ancestor (MRCA) of lineage BD-1
136  was isolated in 1987 (95% HPD: 1983-1991), and lineage BD-2 in 1997 (95% HPD: 1994-2000),
137  where HPD stands for height posterior density. Strains of BD-1 and BD-2 shared genome
138  sequences of strains isolated since 1981 (95% HPD: 1976-1986). The number of SNPs in strains
139  of the two clades is relative to reference *V. cholerae* N16961, which showed strains of BD-0
140  differed by 107 - 137 SNPs, BD-1 by 123 - 189 SNPs, and BD-2 by 146 - 186 SNPs. An
141  unrooted tree showed SNP diversity among BD-0, BD-1, and BD-2 clades with SNP diversity of
142  BD-2 highest (Fig. 1B). Comparison of isolates in the clades and year of isolation revealed
143  clonal aggregation within the dominant clade and strong temporal signature. Strains of BD-1 and
144  BD-2 were found to be temporally spread but simultaneously isolated during the periods of 2004
145  - 2011, 2012, 2014 - 2016 (Fig. 1C, Table S2 supplemental material). Strains of BD-1 were
146  mainly isolated during 2004-2011 (66.3%, n=65) while strains of BD-2 were isolated during
147  those years in fewer numbers (33.7%, n=33) except 2009 when BD-2 strains were dominant
148  (93.33%, n=14) (see Table S1 in supplementary material). The following years, from 2012 to
149  2017, showed BD-2 strains to be dominant (73.5%, n=72) and BD-1 strains the minority (10.2%,
150  n=10).

151  **Genetic variants associated with the clades**

152  Associations between lineages and the genetic variants was studied using 1298 SNPs and 413
153  indels, identified by aligning raw reads against *V. cholerae* N16961 reference genome. Variant

154   annotation using SnpEff (20) showed that among the 1298 SNPs, there were 337 synonymous,

155   613 nonsynonymous, and 348 variants on intergenic regions (Fig. 2A-C, see Table S2 in the

156   supplemental materials). Moreover, of 413 indels, there were 238 frameshift-variants, 107

157   variants on intergenic regions, and 68 other types of variants (Fig 2D-F, Table S2). Most of the

158   identified SNPs and indels were located in the protein-coding region, many of which function to

159   change the form of a protein. By plotting distribution of SNP types and indel variants for BD-0

160   (n=11), BD-1 (n=76), and BD-2 (n=105), it was observed that strains of the clades accumulated

161   SNPs and indels. Strains of BD-2 accumulated more SNPs and indels, increasing genetic

162   distance from BD-0 and BD-1 (Fig. 1B, Fig. 2) and suggesting evolution was occurring when

163   compared with reference *V. cholerae* O1 N16961.

164   Fisher exact test (21) was performed for association analysis between genetic variants and the

165   clades BD-1 and BD-2. Association analysis showed that 140 SNPs and 31 indels had a genome-

166   wide significant association ($p < 6.40 \times 10^{-9}$) with BD-1 and BD-2. Among the 140 SNPs were 25

167   synonymous variants, 53 missense variants, 2 stop gain variants, and 60 variants on intergenic

168   regions (Table S3 and Fig. S2 in supplementary material). It was discovered that 21 SNP

169   missense mutations were present in genes with known functions in more than 80% of BD2

170   strains, resulting in mutant proteins (Table 1). However, there were only seven missense

171   mutations were found in genes with known functions in more than 80% of BD1 strains.

172   Genotype and frequency of 140 significantly associated SNPs, number of SNPs by year of

173   isolation, and root to tip distance, showed significant genetic differences between BD-1 and BD-

174   2 (Fig. 3). The number of core genome SNPs by year of isolation was analyzed to detect

175   temporal SNP accumulation patterns of the clades. The number of core genome SNPs did

176   increase over time for both BD-1 and BD-2 (Fig. 3B). Moreover, root-to-tip regression analysis

177   indicated a steady increase in SNP divergence among the strains of the two clades over time

178   (Fig. 3C). Miami plot for frequency of alternative alleles of the 140 significant SNPs showed

179   BD-2 strains had accumulated more clade specific SNPs, notably in chromosome-2 compared to

180   BD-1 (Fig. 3D).

181   **Relative gene abundance**

182   Pangenome analysis was done using Roary to investigate differences in core and pan genes

183   among the strains of BD-0, BD-1, and BD-2. Roary classified the identified functional genes into

184　　four categories: (i) core genes, present in 99-100% of the strains; (ii) softcore genes, present in

185　　95-99% of the strains; (iii) shell genes, present in 15-95% of the strains; and (iv) cloud genes,

186　　present in less than 15% of the strains (22). Pangenome analysis revealed significant differences

187　　in overall gene composition among the clades (Fig 4A). According to the definition of core genes

188　　in pangenome analysis, the number of core genes largely varied among BD-0, BD-1, and BD-2

189　　(see Table S4 in supplementary material). Similarly, the number of soft-core genes was also

190　　varied. BD-0 is a group of close relatives with a larger genetic distance relative to BD-1 and BD-

191　　2. All BD-0 strains and more than 95% of the BD-1 and BD-2 strains had 1102 common genes

192　　(see Table S5A in supplementary material) most having known function. About 10% of BD-2

193　　strains had 44 unique genes of which six encoding crucial proteins of known function were

194　　found in more than 90% of the BD-2 strains. Those genes are: tetracycline repressor protein

195　　(**tetR**), tetracycline resistance protein (**tetA**), type-I restriction enzyme EcoKI M protein (**hsdM**),

196　　type-I restriction enzyme EcoR124II R protein (**hsdR**), Mrr restriction system protein (**mrr**), and

197　　5-methylcytosine-specific restriction enzyme B (**mcrB**) (see Table S5B in supplementary

198　　material). In addition, methyl-accepting chemotaxis protein (**CtpH**) and group_10030 virulence

199　　proteins were exclusively found in 60% and 65% of BD-2 strains, respectively. By contrast,

200　　about 5-15% of the BD-1 strains carried 19 genes that were unique for them (see Table S5C in

201　　supplementary material). Three genes common to all BD-0 strains were not detected in BD-2 and

202　　were present only in 1-2 of the BD-1 strains.

203　　Next, we conducted Pan-GWAS to identify clade-specific genes by considering gene presence

204　　and absence as the explanatory variable and defined lineage groups as the response variable. A

205　　total of 92 genes were significantly (p-value $< 4.98 \times 10^{-6}$) associated with BD-0 and BD-1 (see

206　　Table S6A in supplementary material). Of these, 62 genes were identified in 54-73% of BD-0

207　　but not in BD-1 strains. Of 164 genes associated with BD-0 and BD-2, 46 were found in more

208　　than 73% of BD-2, but not in BD-0 strains (see Table S6B in supplementary material). In

209　　addition, 66 genes were found in more than 45% of BD-0, but not in BD-2 strains. Of 143 genes

210　　associated with BD-1 and BD-2 (see Table S6C in supplementary material), 29 were found in

211　　more than 76% of BD-1, but not in BD-2 strains. Again, 47 genes were found in 22-97% of the

212　　BD-2, but not in BD-1 strains. These results provide evidence that strains of BD-1 and BD-2

213　　diverged and evolved as two lineages by accumulating genes, after originating from common

214　　ancestor BD-0.

**Pathogenicity islands and phage inducible chromosomal island like elements**

*V. cholerae* strains included in this study were further examined by targeting the pandemic and pathogenicity islands namely VSP-1, VSP-II, VPI-1, and VPI-2, including the phage inducible chromosomal island like elements (PLE). Based on the extent of detected regions compared to *V. cholerae* N16961, five variants of VSP-II (variants 1-5 of the wild type) as reported in our recent study (16), and one variant of VPI-1 (variant 1 of the wild type) were observed (Fig 5). *V. cholerae* El Tor strains differed in type of VSP-II and VPI-1 variants. BD-0 had wild type of VSP-II, as in reference El Tor N16961 strain. Most BD-1 strains (except two) had variant-4 VSP-II, with partial deletion in VC_495 and complete deletion in VC_496 to VC_512, and BD-2 strains carried three VSP-II variants of which ca. 73% had variant-2 VSP-II with partial ORF VC_495 deletion, and complete VC_496 to VC_500 deletion, which appeared consistent with our prior study (16). BD-0 and BD-1 harbored wild type of VPI-1, whereas most of the BD-2 strains (102 of 105 strains) had variant VPI-1 with complete deletion of VC_819 to VC_820 ORFs; and partial deletion in VC_821. All BD-0 strains, and 66 of 76 BD-1 strains lacked PLE (see Tables S1 and S7 in supplementary material), while PLE2 was found in ten BD-1 strains isolated in 2007 possessing the *ctxB*1 genotype and one in 2005. Interestingly, most of the BD-2 strains (83 of 103) carried PLE1, but the rest lacked PLE. Thus, BD-2 lineage strains associated with recent Bangladesh endemic cholera are variant-3 VSP-II, variant VPI-1, and the majority possesses PLE1.

**Variations in SXT/R391 and important genes**

Although differences in SXT/R391, *ctxB*, *gyrA*, *rtxA*, and *parC* across two lineages (BD-1, analogue of lineage-2; BD-2, analogue of lineage-1) were investigated in our recent study (16), these important genetic elements were rechecked to draw overall conclusions for all strains included in this investigation. Moreover, variation in ToxR binding repeats were checked across strains of different lineages. Integrative and conjugative elements (ICEs) were targeted from whole-genome sequences by aligning raw reads or contigs with five publicly available sequences of the ICE element (Accession ID: GQ463140.1, GQ463141.1, GQ463142.1, MK165649.1, and MK165650.1). Nucleotide blast was used to match extracted sequences with ICE element sequences and typed based on highest bit score. Four strains of BD-0 blast search yielded high

244  bit scores when aligned with ICE$^{GEN}$ (MK165650.1), ICE*Vch*Ind5 (GQ463142.1), or

245  ICE*Vch*Ban5 (GQ463140.1). Bit scores were highest for the other BD-0 strains when aligned

246  with ICE$^{TET}$ (Accession ID: MK165649.1), which has genomic characteristics similar to

247  ICEVchVhn2255 (Accession ID: KT151660). For all BD-1 strain bit scores were high when

248  aligned with ICE$^{GEN}$, ICEVchInd5, or ICEVchBan5, and for BD-2 strains bit scores were highest

249  when aligned with ICE$^{TET}$, which is consistent with our previous results. All BD-1 and BD-2

250  strains contained mutant *gyrA* with an amino acid alteration Ser83Ile, whereas 99 (94.28 percent)

251  of the 105 BD-2 strains exhibited Asp660Glu, which was not present in BD-1 or BD-0, also

252  supporting our previous findings.

253  *V. cholerae* O1 El Tor strains in this study were CTX positive, and each carried a single copy of

254  CTXΦ with a particular *ctxB* genotype. Three variants, *ctxB1* (classical genotype), *ctxB3* (typical

255  El Tor genotype), and *ctxB*7 (Haitian variant), of the cholera toxin gene were detected and found

256  associated with the clades (Fig. 1A). Similar to previous findings, all BD-2 strains had *ctxB1*

257  genotype, majority of BD-1 strains had *ctxB7* genotype, and all but two BD-1 strains possessed

258  *rtxA* that differed from El Tor reference N16961 by a single SNP at position 13602 of 1563748

259  bp (NCBI Accession ID: NC 002505.1), corresponding to *rtxA* allele 4 (23). However, in this

260  study it was observed that early BD-1 strains had the *ctxB1* genotype, and over time gained the

261  *ctxB7* genotype.

262  A prior study showed that, Kolkata strains had four heptad repeats (TTTTGAT), whereas

263  Haitian strains had five heptad repeats (24). All BD-0 strains had four heptad repeats (Table S1),

264  while most BD-1 strains (93.4%; n=71) had four repeats, and only 5.3% (n=4) strains had five

265  repeats. As a result, the majority of BD-1 strains with *ctxB*7 genotypes differed from Haitian

266  strains in ToxR binding repeats. BD-2 strains had more diversity in ToxR binding repeats with

267  59.0% (n=62) carrying heptad repeats, 24.8% (n=26) five repeats, and 16.2% (n=17) three

268  repeats.

**Discussion**

270  *Vibrio cholerae* biotype El Tor, the causative agent of the 7$^{th}$ cholera pandemic has increased

271  transmissibility and is more virulent than classical biotype (14, 15). The 7th pandemic strains of

272  cholera circulating in Asia comprises two El Tor clades, one dominant in Bangladesh and the

9

273   other in India (16). Genomic analyses that included additional strains and publicly available
274   genome sequences of wave-2 and wave-3 strains (6, 12) provide a detailed view of longitudinally
275   and temporally representative *V. cholerae* clades associated with endemic cholera in Bangladesh
276   over a period of 27 years (1991 – 2017). The results provide new insights potentially
277   interpretable as origin and progression, based on differences in SNPs, indels, and gene
278   acquisition, including antibiotic resistance cassettes in BD-1 and BD-2, the latter having gained
279   ascendency and dominance as the agent of Bangladesh endemic cholera.

280   Results of whole genome sequencing (16), combined with additional genome sequence data for
281   *V. cholerae* El Tor isolates of Bangladesh endemic cholera, allowed identification of two
282   lineages, designated BD-1 and BD-2. The two clades appear to have originated from a common
283   ancestor of paraphyletic group BD-0, as early as 1981 (95% HPD: 1976-1986). According to A.
284   Mutreja et al. (12), seven strains of BD-0 isolated between 1991 and 2000 represent wave-2
285   strains, and only one strain isolated in 1994, wave-3 with a most recent common ancestor
286   (MRCA) for BD-1 and BD-2. The BD-1 and BD-2 clades may belong to wave-3. Although BD-
287   0 consisted of predominantly of wave-2 strains, three sequenced strains isolated in 2012 shared a
288   wave-2-like genetic background (6), suggesting wave-2 strains may have already been present.
289   Almost all wave-3 strains from a previous study (12) grouped with strains belonging to BD-1.
290   Consistent with results of a previous study (16), significant differences were noted between BD-
291   1 and BD-2, which varied in temporal predominance as the causal agent of Bangladesh endemic
292   cholera. Most (n=62; 82 percent) BD-1 strains had been isolated between 2007 and 2012, with
293   predominance during that time. Between 2005 and 2017, 105 strains belonging to BD-2 were
294   reported, with 97 obtained between 2009 and 2017, implying BD-2 association with recent
295   Bangladesh endemic cholera until 2017. Phylodynamic analysis using BEAST (19) revealed
296   strains of BD-1 had been isolated in Bangladesh roughly ten years before BD-2 strains (see, Fig.
297   S1 in supplementary material), and previously identified as Asian lineage -2 and Asian lineage-1,
298   respectively (16).

299   BD-1 and BD-2 strains appear to have advanced by accumulating different SNPs and indels.
300   Fisher exact test (21) identified 140 SNPs and 31 indel differences between BD-1 and BD-2,
301   resulting in gene alleles unique to them (Fig 3). The majority of the SNPs and indels were
302   components of protein coding genes, suggesting a possible crucial role in their adaption in

303   Bangladesh. Regression analysis of the number of SNPs and year of isolation suggested that both

304   clades consistently accumulated SNPs over time, implying evolution in response to

305   environmental selective pressure.

306   Pangenome analysis using Roary (22) provided evidence of gene acquisition by strains of the

307   clades. A recent study of *V. cholerae* O1 strains isolated in Pakistan found evidence of gene

308   acquisition, where the number of core and accessory genes varied among different lineages (25).

309   According to results of the analysis reported here, the number of core and accessory genes varied

310   significantly among strains of BD-0, BD-1, and BD-2 in Bangladesh (Fig. 4A). The Pan-GWAS

311   approach helped identify genes unique for each clade which could be considered contributing to

312   virulence and/or niche adaptation (26).

313   Phage inducible chromosomal island like elements (PLE) protect *V. cholerae* populations from

314   ICP1 infection by acting as an abortive infection system (27). In this study, the observed

315   predominance in BD-2 of PLE1, not found in BD-0 and BD-1, could have provided a selective

316   advantage for the lineage over BD-1, establishing dominance as an etiological agent of endemic

317   cholera in Bangladesh in recent years.

318   Two BD-0 strains carried CTX phage with *ctxB*3, while other strains carried CTX phage with

319   typical *ctxB*1. Strains at the base of BD-1 had CTX with *ctxB*1 isolated before 2007 and

320   comprised multiple clusters. Moreover, CTX phage of all BD-2 strains contained classical *ctxB*1.

321   A mutation in *rtxA* creating a premature stop codon disabled toxin function in emerging *V.*

322   *cholerae* El Tor strains bearing *ctxB*1 (24). As in the classical strains, altered El Tor pandemic

323   strains eliminated *rtxA* after acquiring classical *ctxB*. In this study, BD-0 and BD-2 strains

324   contained the wild-type *rtxA* allele 1 (Fig. 3A) described by Dolores and Satchell (23). None

325   contained deletions in *rstB* gene when reads were compared to *V. cholerae* N16961 reference

326   genome, indicating *rstB* of Bangladesh *V. cholerae* O1 El Tor isolates does not resemble that of

327   the Haitian outbreak isolates that have been analyzed.

328   ToxR is a global transcriptional regulator of virulence gene expression and this repeated

329   sequence is required for ToxR binding and activation of the *ctxAB* promoter. The ToxR-binding

330   site is located immediately upstream of *ctxAB* and the affinity of ToxR binding is influenced by

331   the repeat sequences (28). The presence of an increased number of ToxR binding repeats located

11

332 between *zot* and *ctxA* has been hypothesized to correlate with a severe form of cholera (28). In
333 this study, variation was detected in the number of ToxR binding repeats (TTTTGAT) among
334 sequences of the *V. cholerae* El Tor isolates. All BD-0 strains had four heptad repeats observed
335 in 93.4% of BD-1 and 59% of BD-2 strains. For BD-2 strains, however, greater variation was
336 observed in ToxR binding repeats as ca. 24.8% (n=26) of BD-2 strains contained five heptad
337 repeats, whereas 16.2% (n=17) had three heptad repeats, suggesting robustness of the clade.

338 Targets of quinolones are type II topoisomerases of DNA gyrase, a heterotetramer composed of
339 two A and two B subunits, encoded by *gyrA* and *gyrB* genes respectively (29). It was observed
340 that all BD-1 and BD-2 strains had a common mutation Ser83 to Ile in *gyrA*, while 94.29%
341 (99/105) BD-2 had an additional mutation Asp660 to Glu. Furthermore, 87% (66/76) of BD-1
342 strains exhibited a mutation Ser85 to Leu *parC*, whereas all BD-2 strains (105/105) had this
343 mutation. In Haitian *V. cholerae* strains, *gyrA* and *parC* genes had two point mutations: Ser83 to
344 Ile in *gyrA* and Ser85 to Leu in *parC*. Both are linked to quinolone resistance in *V. cholerae*
345 strains associated with recent cholera outbreaks in India, Nigeria, and Cameroon (30).

346 SXT/R391 family ICEs are transferable elements associated with antimicrobial resistance in *V.*
347 *cholerae* (31). The SXT-ICE regions of the isolates included in this study, were compared with
348 five sequences of the elements to the type SXT/R391 family ICEs belonging to strains associated
349 with cholera (*V. cholerae* O1 and O139) (9, 32). Four BD-0 strains exhibited ICE elements
350 similar to ICE$^{GEN}$, ICEVchInd5 or ICEVchBan5, whereas the rest had ICE elements similar to
351 ICE$^{TET}$. Interestingly, ICE elements of BD-1 strains included ICE$^{GEN}$, ICEVchInd5 or
352 ICEVchBan5-like ICE elements, whereas BD-2 strains differed completely from the others, with
353 only ICE$^{TET}$-like ICE elements.

354 The results of the study reported here included BD-1 and BD-2 isolated during the Bangladesh
355 endemic cholera of 2004 onwards and that, while existing together, with each subsequent year
356 they exhibited different dominance. BD-2 diverged, while retaining the ability to produce
357 multifunctional-autoprocessing repeats-in-toxin (MARTX) and acquiring SXT element ICE$^{TET}$
358 containing tetracycline resistance genes. This observation hints at a selective advantage of BD-2
359 strains over BD-1 strains for robustness. It is evident from results of the analyses that BD-1 and
360 BD-2 differ significantly, owing to gene composition and SNPs and may have evolved
361 independently due to selection pressures. The use of antibiotics, including tetracycline, can exert

12

362     selection pressure in evolution (16, 33), while strains stopping to produce MARTX along with

363     other variations in the genome might provide a selective advantage. According to suggestions

364     from studies of the dynamics of *V. cholerae*, immunocompetence of the host against *V. cholerae*

365     strains may contribute to the dynamics of *V. cholerae*, hence produce an effect from interaction

366     with humans in selection and cannot be ruled out (34).

367     Cholera globally is influenced by thriving populations of *V. cholerae* occurring naturally in the

368     Ganges Delta of Bay of Bengal (GDBB) (1, 2, 5, 14). Overall results presented here suggest

369     means of emergence and progression of the two clades in evolution from a progenitor *V.*

370     *cholerae* El Tor initiating the seventh pandemic in Asia (5) and reflecting short-term evolution

371     of *V. cholerae* El Tor associated with Bangladesh endemic cholera in the GDBB (14, 31). BD-2

372     is concluded to have emerged relatively recently and evolves by acquiring SNPs over time. Also,

373     BD-2 strains showed diversity in indels, possessing SXT/R391 family ICE-elements, PLE1, *tetR*,

374     and several other important genetic elements, and predominantly associated with recent

375     Bangladesh endemic cholera. As is apparent from our results, BD-1 appears to be an analogue of

376     a previously reported lineage 2 from Asia, the major causative agent of cholera in India, Yemen,

377     and Haiti (16). In contrast, BD-2 strains of the present study appear to be an analogue of Asian

378     lineage 1, which successfully outcompeted BD-1 (Asian lineage 2) and established

379     predominance as an etiological agent of cholera in an historical hotspot of the disease,

380     Bangladesh. It can be concluded that this is a reflection of robustness of BD-2 as an epidemic

381     clone emerging locally with potential to transmit globally, and underscoring the need to track the

382     two successful *V. cholerae* El Tor clades.

383     **Materials and Methods**

384     **Bacterial isolates**

385     A total of 119 *V. cholerae* O1 strains from the icddr,b collection of strains isolated in Bangladesh

386     between 2004 and 2017 (see Table S1 in supplementary material) were sequenced. Paired-end

387     Illumina short reads for the isolated strains were generated (150 bp, 150 bp) using MiSeq or

388     Hiseq 2500 sequencer as described in our recent study (16). Publicly available paired-end raw

389     reads of 17 strains isolated in Bangladesh between 1991 and 2007 (see study flow chart Fig. S3

13

390  in supplementary material) and 56 strains from our recent study (16) were included in the
391  analysis.

**Genome assembly and gene annotation**

393  An ultra-fast FASTQ preprocessor implemented in FASTP (35), was used to inspect raw paired-
394  end reads and filter bad ligation or adapter parts. De novo genome assembly implemented in
395  VelvetOptimizer (36) was used to build contigs by optimizing the parameter N50, a metric for
396  assessing contiguity of an assembly. The bacterial genome annotation tool, Prokka (37), was
397  used for whole-genome gene annotation. ResFinder (38) was used to find the antimicrobial
398  resistant gene profiles for all of the strains.

**SNP identification and phylogenetic analysis**

400  Bowtie2 (39) was used to align high-quality reads with reference genome sequence of *V.*
401  *cholerae* N16961 El Tor (NCBI Accession ID: NC_002505.1 and NC_002506.1) for variant
402  calling. Samtools (40) and Bcftools (41) were used to call genome variants. A maximum-
403  likelihood phylogeny was inferred on an alignment of concatenated SNPs evenly distributed
404  across non-repetitive, non-recombinant core genome using IQ-TREE v1.6.1 (42). Trees were
405  visualized in FigTree v1.4.3 (http://tree.bio.ed.ac.uk/ software/figtree/) or Interactive Tree of Life
406  online tool (43).

407

**Bayesian phylogenetic inference**

409  The Bayesian Evolutionary Analysis Sampling Trees (BEAST) v.2.4.4 software package (19)
410  was used for temporal analysis to estimate divergence date of *V. cholerae* O1 isolates in
411  Bangladesh. The date of isolation of each strain was used as tip data. A random clock model was
412  implemented using Markov Chain Monte Carlo (MCMC) chains run for 100 million generations
413  with 10% burn-in and sampled every 1000 generations. A GTR nucleotide substitution model
414  was used. Tree data were summarized using TreeAnotator, a tool of BEAST software package,
415  to generate the maximum clade credibility tree.

**Pangenome analysis**

417    A pan-genome was constructed using Roary (22) from annotated assemblies of the sample set
418    with percentage protein identity of 95%. The protein sequences were first extracted and
419    iteratively pre-clustered with cd-hit (version 4.6) down to 98% identity. An all against all blast
420    (version 2.2.31) was performed on the remaining non-clustered sequences and a single
421    representative sequence from each cd-hit cluster was selected. The data were used by MCL (44)
422    (version 11–294) to cluster the sequences. The preclusters and MCL clusters were merged and
423    paralogs split by inspecting the conserved gene neighborhood around each sequence (5 genes on
424    either side). Each sequence for each cluster was independently aligned using PRANK (45)
425    (version 0.140603) and combined to form a multi-FASTA alignment of the core genes.
426    Sequences of SXT elements were compared with ICE$^{GEN}$ and ICE$^{TET}$ using BRIG 0.95 with 70%
427    BLAST identity (46).

428

438    **Reference**

439
440    1.    Reidl J, Klose KE. 2002. Vibrio cholerae and cholera: out of the water and into the host. FEMS
441          Microbiol Rev 26:125-39.
442    2.    Albert MJ, Siddique A, Islam M, Faruque A, Ansaruzzaman M, Faruque S, Sack RB. 1993. Large
443          outbreak of clinical cholera due to Vibrio cholerae non-01 in Bangladesh. The Lancet 341:704.
444    3.    Ramamurthy T, Das B, Chakraborty S, Mukhopadhyay AK, Sack DA. 2019. Diagnostic techniques
445          for rapid detection of Vibrio cholerae O1/O139. Vaccine doi:10.1016/j.vaccine.2019.07.099.
446    4.    Clemens JD, Nair GB, Ahmed T, Qadri F, Holmgren J. 2017. Cholera. The Lancet 390:1539-1549.
447    5.    Hu D, Liu B, Feng L, Ding P, Guo X, Wang M, Cao B, Reeves PR, Wang L. 2016. Origins of the
448          current seventh cholera pandemic. Proc Natl Acad Sci U S A 113:E7730-E7739.

449  6.   Domman D, Quilici ML, Dorman MJ, Njamkepo E, Mutreja A, Mather AE, Delgado G, Morales-
450       Espinosa R, Grimont PAD, Lizarraga-Partida ML, Bouchier C, Aanensen DM, Kuri-Morales P, Tarr
451       CL, Dougan G, Parkhill J, Campos J, Cravioto A, Weill FX, Thomson NR. 2017. Integrated view of
452       Vibrio cholerae in the Americas. Science 358:789-793.
453  7.   Kim EJ, Lee D, Moon SH, Lee CH, Kim SJ, Lee JH, Kim JO, Song M, Das B, Clemens JD, Pape JW,
454       Nair GB, Kim DW. 2014. Molecular insights into the evolutionary pathway of Vibrio cholerae O1
455       atypical El Tor variants. PLoS Pathog 10:e1004384.
456  8.   Rashid MU, Rashed SM, Islam T, Johura FT, Watanabe H, Ohnishi M, Alam M. 2016. CtxB1
457       outcompetes CtxB7 in Vibrio cholerae O1, Bangladesh. J Med Microbiol 65:101-103.
458  9.   Wozniak RA, Fouts DE, Spagnoletti M, Colombo MM, Ceccarelli D, Garriss G, Dery C, Burrus V,
459       Waldor MK. 2009. Comparative ICE genomics: insights into the evolution of the SXT/R391 family
460       of ICEs. PLoS Genet 5:e1000786.
461  10.  Faruque SM, Mekalanos JJ. 2003. Pathogenicity islands and phages in Vibrio cholerae evolution.
462       Trends Microbiol 11:505-10.
463  11.  Murphy RA, Boyd EF. 2008. Three pathogenicity islands of Vibrio cholerae can excise from the
464       chromosome and form circular intermediates. J Bacteriol 190:636-47.
465  12.  Mutreja A, Kim DW, Thomson NR, Connor TR, Lee JH, Kariuki S, Croucher NJ, Choi SY, Harris SR,
466       Lebens M, Niyogi SK, Kim EJ, Ramamurthy T, Chun J, Wood JL, Clemens JD, Czerkinsky C, Nair GB,
467       Holmgren J, Parkhill J, Dougan G. 2011. Evidence for several waves of global transmission in the
468       seventh cholera pandemic. Nature 477:462-5.
469  13.  Islam MT, Clemens JD, Qadri F. 2018. Cholera Control and Prevention in Bangladesh: An
470       Evaluation of the Situation and Solutions. J Infect Dis 218:S171-S172.
471  14.  Nair GB, Qadri F, Holmgren J, Svennerholm AM, Safa A, Bhuiyan NA, Ahmad QS, Faruque SM,
472       Faruque AS, Takeda Y, Sack DA. 2006. Cholera due to altered El Tor strains of Vibrio cholerae O1
473       in Bangladesh. J Clin Microbiol 44:4211-3.
474  15.  Taneja N, Mishra A, Sangar G, Singh G, Sharma M. 2009. Outbreaks caused by new variants of
475       Vibrio cholerae O1 El Tor, India. Emerg Infect Dis 15:352-4.
476  16.  Morita D, Morita M, Alam M, Mukhopadhyay AK, Johura F-t, Sultana M, Monira S, Ahmed N,
477       Chowdhury G, Dutta S. 2020. Whole-Genome Analysis of Clinical Vibrio cholerae O1 in Kolkata,
478       India, and Dhaka, Bangladesh, Reveals Two Lineages of Circulating Strains, Indicating Variation in
479       Genomic Attributes. Mbio 11.
480  17.  Leinonen R, Akhtar R, Birney E, Bower L, Cerdeno-Tarraga A, Cheng Y, Cleland I, Faruque N,
481       Goodgame N, Gibson R, Hoad G, Jang M, Pakseresht N, Plaister S, Radhakrishnan R, Reddy K,
482       Sobhany S, Ten Hoopen P, Vaughan R, Zalunin V, Cochrane G. 2011. The European Nucleotide
483       Archive. Nucleic Acids Res 39:D28-31.
484  18.  Baddam R, Sarker N, Ahmed D, Mazumder R, Abdullah A, Morshed R, Hussain A, Begum S,
485       Shahrin L, Khan AI, Islam MS, Ahmed T, Alam M, Clemens JD, Ahmed N. 2020. Genome Dynamics
486       of Vibrio cholerae Isolates Linked to Seasonal Outbreaks of Cholera in Dhaka, Bangladesh. mBio
487       11.
488  19.  Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. 2018. Bayesian
489       phylogenetic and phylodynamic data integration using BEAST 1.10. Virus Evol 4:vey016.
490  20.  Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012. A
491       program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:
492       SNPs in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3. Fly (Austin) 6:80-92.
493  21.  Raymond M, Rousset F. 1995. An exact test for population differentiation. Evolution:1280-1283.
494  22.  Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M, Falush D, Keane JA,
495       Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome analysis. Bioinformatics
496       31:3691-3.

23. Dolores J, Satchell KJ. 2013. Analysis of Vibrio cholerae genome sequences reveals unique rtxA variants in environmental strains and an rtxA-null mutation in recent altered El Tor isolates. mBio 4:e00624.

24. Ghosh P, Naha A, Pazhani GP, Ramamurthy T, Mukhopadhyay AK. 2014. Genetic traits of Vibrio cholerae O1 Haitian isolates that are absent in contemporary strains from Kolkata, India. PLoS One 9:e112973.

25. Zeb S, Gulfam SM, Bokhari H. 2020. Comparative core/pan genome analysis of Vibrio cholerae isolates from Pakistan. Infect Genet Evol 82:104316.

26. Gori A, Harrison OB, Mlia E, Nishihara Y, Chan JM, Msefula J, Mallewa M, Dube Q, Swarthout TD, Nobbs AH, Maiden MCJ, French N, Heyderman RS. 2020. Pan-GWAS of Streptococcus agalactiae Highlights Lineage-Specific Genes Associated with Virulence and Niche Adaptation. mBio 11.

27. Hays SG, Seed KD. 2020. Dominant Vibrio cholerae phage exhibits lysis inhibition sensitive to disruption by a defensive phage satellite. eLife 9:e53200.

28. Pfau JD, Taylor RK. 1996. Genetic footprint on the ToxR-binding site in the promoter for cholera toxin. Mol Microbiol 20:213-22.

29. Hooper DC. 1998. Clinical applications of quinolones. Biochim Biophys Acta 1400:45-61.

30. Hasan NA, Choi SY, Eppinger M, Clark PW, Chen A, Alam M, Haley BJ, Taviani E, Hine E, Su Q, Tallon LJ, Prosper JB, Furth K, Hoq MM, Li H, Fraser-Liggett CM, Cravioto A, Huq A, Ravel J, Cebula TA, Colwell RR. 2012. Genomic diversity of 2010 Haitian cholera outbreak strains. Proc Natl Acad Sci U S A 109:E2010-7.

31. Weill FX, Domman D, Njamkepo E, Almesbahi AA, Naji M, Nasher SS, Rakesh A, Assiri AM, Sharma NC, Kariuki S, Pourshafie MR, Rauzier J, Abubakar A, Carter JY, Wamala JF, Seguin C, Bouchier C, Malliavin T, Bakhshi B, Abulmaali HHN, Kumar D, Njoroge SM, Malik MR, Kiiru J, Luquero FJ, Azman AS, Ramamurthy T, Thomson NR, Quilici ML. 2019. Genomic insights into the 2016-2017 cholera epidemic in Yemen. Nature 565:230-233.

32. Sarkar A, Morita D, Ghosh A, Chowdhury G, Mukhopadhyay AK, Okamoto K, Ramamurthy T. 2019. Altered Integrative and Conjugative Elements (ICEs) in Recent Vibrio cholerae O1 Isolated From Cholera Cases, Kolkata, India. Front Microbiol 10:2072.

33. Tello A, Austin B, Telfer TC. 2012. Selective pressure of antibiotic pollution on bacteria of importance to public health. Environ Health Perspect 120:1100-6.

34. Chakraborty S, Mukhopadhyay AK, Bhadra RK, Ghosh AN, Mitra R, Shimada T, Yamasaki S, Faruque SM, Takeda Y, Colwell RR, Nair GB. 2000. Virulence genes in environmental strains of Vibrio cholerae. Appl Environ Microbiol 66:4022-8.

35. Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor. Bioinformatics 34:i884-i890.

36. Zerbino DR, Birney E. 2008. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. Genome Res 18:821-9.

37. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics 30:2068-9.

38. Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, Aarestrup FM, Larsen MV. 2012. Identification of acquired antimicrobial resistance genes. J Antimicrob Chemother 67:2640-4.

39. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nat Methods 9:357-9.

40. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25:2078-9.

41. Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. Bioinformatics 27:2987-93.

545    42.    Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective stochastic
546        algorithm for estimating maximum-likelihood phylogenies. Mol Biol Evol 32:268-74.
547    43.    Letunic I, Bork P. 2019. Interactive Tree Of Life (iTOL) v4: recent updates and new developments.
548        Nucleic Acids Res 47:W256-W259.
549    44.    Enright AJ, Van Dongen S, Ouzounis CA. 2002. An efficient algorithm for large-scale detection of
550        protein families. Nucleic Acids Res 30:1575-84.
551    45.    Loytynoja A. 2014. Phylogeny-aware alignment with PRANK. Methods Mol Biol 1079:155-70.
552    46.    Alikhan NF, Petty NK, Ben Zakour NL, Beatson SA. 2011. BLAST Ring Image Generator (BRIG):
553        simple prokaryote genome comparisons. BMC Genomics 12:402.

554

555

556

## Tables

558

559    **Table 1. SNPs resulted unique mutant proteins in BD1 and BD2**

| SNP | REF | ALT | FrqBD1 | FrqBD2 | p-value | Gene | AA change | Product |
|---|---|---|---|---|---|---|---|---|
| S1_2609994 | G | A | 0 | 105 | 5.61E-53 | nudF_1 | Arg109Cys | ADP-ribose pyrophosphatase |
| S2_266019 | A | G | 0 | 105 | 5.61E-53 | ulaA | Ile354Thr | Ascorbate-specific permease IIC component UlaA |
| S2_1024884 | G | A | 0 | 105 | 5.61E-53 | putA | Ala600Val | Bifunctional protein PutA |
| S2_989172 | C | T | 0 | 105 | 5.61E-53 | yecS | Pro191Ser | YecS |
| S1_798976 | T | C | 0 | 105 | 5.61E-53 | suhB | Glu217Gly | Inositol-1-monophosphatase |
| S1_994229 | G | A | 0 | 105 | 5.61E-53 | stcE_2 | Gly201Asp | Metalloprotease StcE precursor |
| S2_921045 | A | C | 0 | 105 | 5.61E-53 | ctpH_6 | Ile161Ser | Methyl-accepting chemotaxis protein CtpH |
| S1_1622584 | G | A | 0 | 105 | 5.61E-53 | cobB | Pro50Leu | NAD-dependent protein deacylase |
| S2_773493 | T | A | 0 | 105 | 5.61E-53 | phhA | Gln19Leu | Phenylalanine-4-hydroxylase |
| S1_681574 | G | T | 0 | 105 | 5.61E-53 | glmM | Arg196Leu | Phosphoglucosamine mutase |
| S2_161094 | T | G | 0 | 105 | 5.61E-53 | siaT_5 | Ser241Ala | Sialic acid TRAP transporter permease protein SiaT |
| S1_1452755 | T | C | 0 | 105 | 5.61E-53 | cysG_1 | Val38Ala | Siroheme synthase |
| S1_2731709 | G | A | 0 | 105 | 5.61E-53 | tamA | Thr266Ile | Translocation and assembly module TamA precursor |
| S1_545919 | T | G | 0 | 104 | 4.32E-51 | pctB_1 | Leu249Trp | Methyl-accepting chemotaxis protein PctB |
| S1_2814292 | T | C | 0 | 102 | 4.43E-48 | argG | Thr283Ala | Argininosuccinate synthase |
| S1_1332186 | T | G | 0 | 99 | 1.96E-44 | gyrA | Asp660Glu | DNA gyrase subunit A |
| S1_149686 | G | T | 0 | 99 | 1.96E-44 | murI | Ala137Ser | Glutamate racemase |
| S2_562858 | A | T | 0 | 99 | 1.96E-44 | VCA0627 | Thr6Ser | rRNA methylase |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| S1_628646 | C | T | 0 | 85 | 1.32E-32 | hrpB_1 | Ala782Val | ATP-dependent RNA helicase HrpB |
| S1_673206 | A | G | 0 | 85 | 1.32E-32 | tyrS_2 | Thr393Ala | Tyrosine--tRNA ligase |
| S1_2357516 | G | A | 0 | 79 | 7.24E-29 | angR | Leu227Phe | Anguibactin system regulator |
| S1_2483236 | G | A | 66 | 0 | 4.18E-39 | lysX | Ala150Thr | Alpha-aminoadipate--LysW ligase LysX |
| S1_1682925 | C | T | 67 | 0 | 3.63E-40 | appC | Ala226Thr | Cytochrome bd-II ubiquinol oxidase subunit 1 |
| S1_368119 | T | C | 67 | 0 | 3.63E-40 | mutL | Cys350Arg | DNA mismatch repair protein MutL |
| S1_1359179 | G | A | 67 | 0 | 3.63E-40 | licH | Ala56Thr | putative 6-phospho-beta-glucosidase |
| S1_1060408 | C | T | 71 | 0 | 6.86E-45 | nagA_1 | Asp150Asn | N-acetylglucosamine-6-phosphate deacetylase |
| S1_276112 | G | A | 76 | 0 | 5.61E-53 | mak | Gly116Arg | Fructokinase |
| S1_1782501 | G | A | 76 | 0 | 5.61E-53 | cph2_4 | Leu79Phe | Phytochrome-like protein cph2 |

560  Here, SNP refers the SNPs which had alternative allele uniquely found in more than 80% of BD1 or BD-2 strains, located within
561  proteins of known functions, and alter amino acid. SNPs were named according its chromosomal position. For example,
562  "S1_2609994" is an SNP/indel site, where "S" stands for site and "2609994" stands for the site's base pair location. Reference
563  allele = REF, alternative allele = ALT, AA change = amino acid change. Freq_BD1 is frequency of alternative allele in BD1 and
564  Freq_BD2 is frequency of alternative allele in BD2. Note that, frequencies of alternative alleles of the SNPs are zero for BD-0. P-
565  value is the p-value of Fisher exact test.

566

## Figure Legends

568  **FIG 1** Phylogenetic analyses of strains showing respective genomic features and year of
569  isolation. (A) Maximum likelihood phylogenetic tree generated from whole genome SNPs and
570  number of isolated *V. cholerae* O1 El Tor strains belonging to lineages BD-0, BD-1, and BD-2
571  rooted from out-group reference strain *Vibrio cholerae* N16961. Rings show features of the
572  isolates according to color scheme provided on the left. Tree branches are colored blue, green,
573  and red defining lineages BD-0, BD-1, and BD-2, respectively; (B) Unrooted tree showing
574  independent evolution of BD-1 and BD-2 strains with the number of core genome SNPs of
575  strains in the lineages compared to the N16961 reference strain; and (C) Percentage of isolates
576  per year for the three lineages. Size of the circles indicates percentage of strains belonging to
577  lineages according to the scheme shown.

578  **FIG 2** Box plots of SNPs distribution and indel type in each of three lineage groups. (A)
579  Distribution of 337 synonymous SNP variants. This figure shows that strains of BD-2 lineage
580  accumulated more synonymous SNP variants compared to BD-0 and BD-1 lineages. Notably,
581  synonymous SNP variants do not change the form of protein. (B) Distribution of 613

582 nonsynonymous SNP variants. These nonsynonymous SNP variants include 570 missense

583 variants, 38 stop gained variants, 2 splice-region-variants and stop-retained-variants, 2 stop-lost

584 and splice-region-variants, 1 initiator codon variant. (C) Distribution of 348

585 upstream/downstream SNP variants. (D) Distribution of 238 frameshift indel variants. (E)

586 Distribution of 107 upstream/downstream indel variants. (F) Distribution of 68 indel variants,

587 including 13 conservative-inframe-insertions, 14 disruptive-inframe-insertions, 11 frameshift-

588 variant and stop-gained, 10 disruptive-inframe-deletions, 10 conservative-inframe-deletions, 1

589 stop-gained and disruptive-inframe-deletions, 2 feature-elongations, 1 frameshift-variant and

590 stop-lost and splice-region-variant, 1 stop-gained and disruptive-inframe-insertion, 2 frameshift-

591 variant and splice-region-variant, 2 frameshift-variant and start-lost, 1 stop-gained and

592 conservative-inframe-insertion.

593 **FIG 3** SNP analysis of genetic diversity. (A) Phylogenetic tree map of the strains and heat map

594 for genotypes of 140 SNPs significantly associated with different lineages. Colors used delineate

595 four different nucleotides where white represents the missing genotype. Heatmap shows clear

596 differences in the lineages. (B) Number of core genome SNPs referencing the year of isolation.

597 The figure shows steady accumulation of SNPs of different lineage strains over time. (C)

598 Regression analysis of root-to-tip distance for strains of the lineages. This figure shows diversity

599 of strains of different lineages. (D) Miami plot of alternative allele frequencies of SNPs for the

600 dominant lineages BD-1 and BD-2. This figure shows the clear difference in SNP accumulation

601 by the two dominant lineages BD-1 and BD-2.

602 **FIG 4** Pangenome analysis showing differences in abundance of gene clusters among the

603 lineages. (A) Relative gene abundance of lineages identified by Roary. Features of the sequences

604 are shown with bars and details for features listed in Table S1. (B) BLAST coverage of SXT

605 regions of BD-1 isolates compared with ICE-GEN. Rings represent sequentially outwards

606 following Table S1. Outermost ring shows the different genes of ICE-GEN. (C) BLAST

607 coverage of SXT regions of BD-2 isolates compared with ICE-TET. The rings represent strains

608 of BD-2 sequentially outwards following Table S1. The outermost ring shows different genes of

609 ICE-TET.

610 **FIG 5** Schematic diagram of VSP-II. Schematic alignment view of VSP-II regions for the

611 isolates. Direction of gene transcription is indicated by arrows and gene shadows represent

20

612 functional annotation. Six types were identified with all BD-0 strains wild-type VSP-II. Two

613 major types, var-2 and var-3, observed for most BD-2 strains and one major type var-4 for most

614 BD-1 strains.

## Supplementary Information

**FIG S1 Bayesian phylogenetic analysis of *V. cholerae* O1.** Node ages obtained from BEAST

analysis. Tree visualized using FigTree v1.4.4. Colors of clades reference the lineage.

**FIG S2 Manhattan plots of *p*-values for association studies of SNPs and BD-1 and BD-2**

**lineages.** Blue represent suggested significant and red indicates high significance. Association

analysis reveals 140 SNP difference between BD-1 and BD-2 lineages.

**FIG S3 Study flow chart.** Data curation and analyses steps are given in the flow chart.

**Table S1. Genetic characteristics of strains included in the study.** Lineage refers to

genetically homogeneous groups of strains. Legends are strain ID, year of isolation, SXT/ICE

elements, acquired antibiotic resistance profile, gyrA allele, number of ToxR binding repeats,

ctxB allele, and PLE are tabulated.

**Table S2. Number of strains belonging to the different lineages.** Here, N_BD-0 = number of

strains belonging to BD-0; N_BD-1 = number of strains belonging to BD-1; and N_BD2 =

number of strains belonging to BD-2.

**Table S3. Fisher exact test identifying significantly associated 140 SNPs and 31 indels of the**

**two dominant lineages, BD-1 and BD-2.** SNP/indel indicates significant SNP/indels identified

by Fisher exact test. SNPs and indels named according to chromosomal position. For example,

"S1 1905668" is an SNP/indel site, where "S" stands for site and "1905668" stands for location

for site base pair. Reference allele is Ref, where as the alternative alleles are Alt1 and Alt2. P-

value is Fisher exact test value. Variant type indicates SNP and indel type.

**Table S4. Roary pangenome analysis showing gene compositions differences by lineage.**
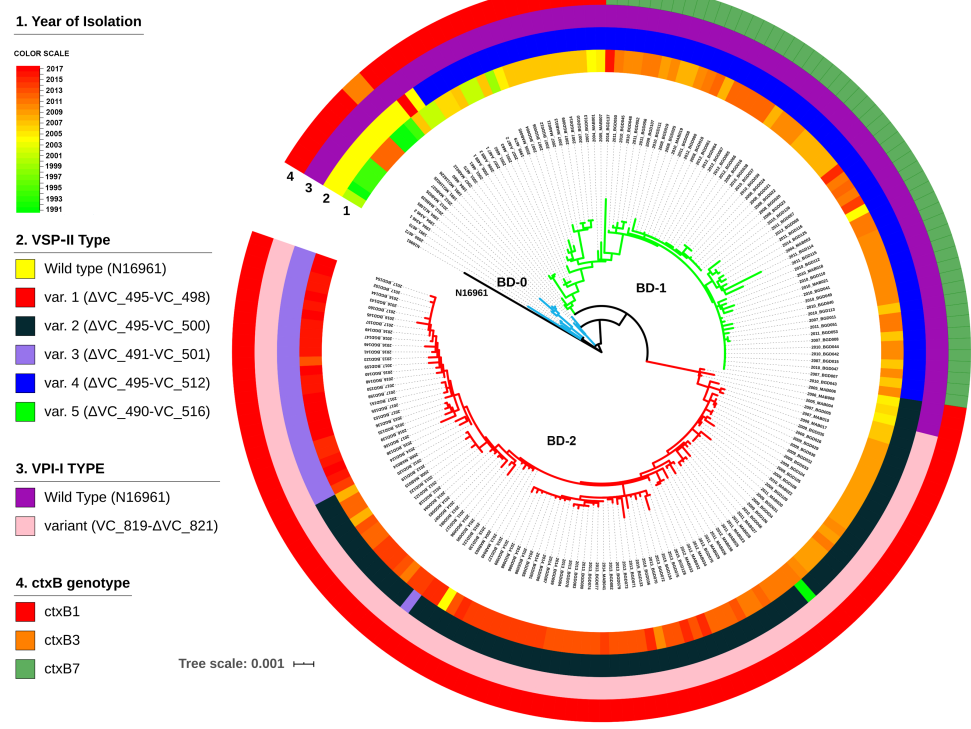
636    Gene cluster refers to group of genes clustered based on existence in the strains of different

637    lineages. Code for number of genes in lineage BD-0, BD-1, and BD-2 is N_BD-0, N_BD-1, and

638    N_BD-2, respectively.

639    **Table S5. Common and unique genes of the different lineages.** (A) Genes detected in more

640    than 95% of BD-0, BD-1, and BD-2 strains. (B) List of unique genes in BD-2. (C) List of unique
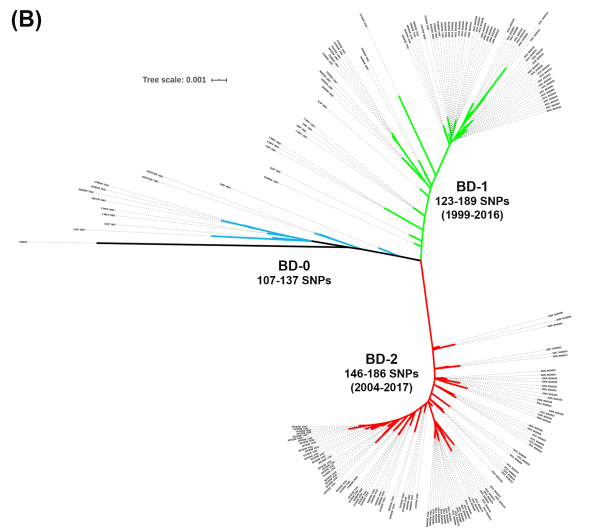
641    genes in BD-1.

642    **Table S6**. **PanGWAS identified lineage associated genes.** (A)  List of genes associated with

643    BD-0 and BD-1. (B) List of genes associated with BD-0 and BD-2. (C) List of genes associated

644    with BD-1 and BD-2.

645    **Table S7. Number of strains with phage inducible chromosomal island like elements (PLE)**.

646    Absence of PLE = PLE(-), and PLE1 and PLE2 are two different types of PLE.