

**Core and accessory genomic traits of *Vibrio Cholerae* O1 drive lineage transmission and disease severity**

Alexandre Maciel-Guerra<sup>1‡</sup>, Kubra Babaarslan<sup>1‡</sup>, Michelle Baker<sup>1‡</sup>, Aura Rahman<sup>2</sup>, Muhammad Maqsud Hossain<sup>1,2</sup>, Abdus Sadique<sup>2</sup>, Jahidul Alam<sup>2</sup>, Salim Uzzaman<sup>3</sup>, Mohammad Ferdous Rahman Sarker<sup>3</sup>, Nasrin Sultana<sup>3</sup>, Ashraful Islam Khan<sup>4</sup>, Yasmin Ara Begum<sup>4</sup>, Mokibul Hassan Afrad<sup>4</sup>, Nicola Senin<sup>5</sup>, Zakir Hossain Habib<sup>3</sup>, Tahmina Shirin<sup>3</sup>, Firdausi Qadri<sup>4</sup>, and Tania Dottorini<sup>1,6\*</sup>

‡ These authors contributed equally

\*Corresponding author (tania.dottorini@nottingham.ac.uk)

The authors declare no competing interests.

<sup>1</sup>School of Veterinary Medicine and Science, University of Nottingham, College Road, Sutton Bonington, Loughborough, Leicestershire, UK, LE12 5RD

<sup>2</sup>North South University, Bashundhara, Dhaka 1229, Bangladesh

<sup>3</sup>Institute of Epidemiology, Disease Control and Research (IEDCR), 44, Shaheed Tajuddin Ahmed Sarani Mohakhali, Dhaka 1212, Bangladesh

<sup>4</sup>International Centre for Diarrhoeal Disease Research (icddr,b), 68, Shaheed Tajuddin Ahmed Sarani Mohakhali, Dhaka 1212, Bangladesh

<sup>5</sup>Department of Engineering, University of Perugia, 06125 Perugia, Italy

<sup>6</sup>Centre for Smart Food Research, Nottingham Ningbo China Beacons of Excellence Research and Innovation Institute, University of Nottingham Ningbo China, Ningbo, P. R. China, 315100

## 1 **Abstract**

2 In Bangladesh, *Vibrio cholerae* lineages are undergoing genomic evolution, with increased virulence  
3 and spreading ability. However, our understanding of the genomic determinants influencing lineage  
4 transmission and disease severity remains incomplete.

5 Here, we developed a computational framework using machine-learning, genome scale metabolic  
6 modelling (GSSM) and 3D structural analysis, to identify *V. cholerae* genomic traits linked to lineage  
7 transmission and disease severity. We analysed in-patients isolates from six Bangladeshi regions (2015-  
8 2021), and uncovered accessory genes and core SNPs unique to the most recent dominant lineage, with  
9 virulence, motility and bacteriophage resistance functions.

10 We also found a strong correlation between *V. cholerae* genomic traits and disease severity, with some  
11 traits overlapping those driving lineage transmission. GSMM and 3D structure analysis unveiled a  
12 complex interplay between transcription regulation, protein interaction and stability, and metabolic  
13 networks, associated to lifestyle adaptation, intestinal colonization, acid tolerance and symptom  
14 severity. Our findings support advancing therapeutics and targeted interventions to mitigate cholera  
15 spread.

16

## 17 **Introduction**

18 Cholera is an acute diarrheal disease. Worldwide, 1.3 billion people are estimated to be at risk and  
19 approximately 1.3 to 4 million cases occur annually, with 21,000 to 143,000 resulting in death<sup>1,2</sup>. In  
20 Bangladesh alone, where cholera is endemic, an estimated 66 million people are at risk of cholera with  
21 at least 100,000 cases and 4,500 deaths per year<sup>1,3</sup>. Globally the O1 serogroup remains the primary  
22 cause of cholera<sup>1,2</sup>. The O1 serogroup is divided into the main serotypes Ogawa and Inaba, and  
23 subdivided into two biotypes, classical and El Tor (7th pandemic), which are genotypically and  
24 phenotypically distinct<sup>4-6</sup>. *V. cholerae* has shown an extraordinary capacity to undergo genetic and  
25 phenotypic changes over time, giving rise to successive waves of genetically and phenotypically diverse

26 pandemic clones. These variants exhibit increased virulence, pathogenicity, resistance and spreading  
27 capability<sup>7,8</sup>.

28 Recently, distinctive lineages belonging to the 7th pandemic El Tor (7PET) wave-3 have been observed  
29 circulating in Bangladesh<sup>9-11</sup>. The two most prominent circulating lineages identified over the last 20  
30 years are BD-1 and BD-2<sup>9-11</sup>, and more recently BD-1.2, responsible for the latest 2022 massive  
31 outbreak in the country<sup>10</sup>. Genomic analysis revealed variations between BD-1.2 and BD-2 in the *Vibrio*  
32 seventh pandemic island II (VSP-II), *Vibrio* pathogenic island 1 (VPI-1), mobile genetic elements,  
33 phage-inducible chromosomal island-like element (PLE), and SXT-related integrating conjugative  
34 elements (SXT ICE)<sup>10</sup>. Despite the advances of genomic analysis, the complete genomic repertoire and  
35 the mechanisms causing the greater transmission of BD-1.2 remain unknown. Gaps persist in our  
36 knowledge regarding whether coding or non-coding single nucleotide polymorphisms (SNPs), or  
37 accessory genes, drive the evolutionary shifts. It remains unclear whether gene regulation, metabolic or  
38 molecular networks, or folding events play a role. There is even less knowledge about the genomic  
39 determinants responsible for the severity of cholera resulting from these lineages. About 1 in 5 people  
40 with cholera will experience a severe condition owing to a combination of symptoms (primarily  
41 diarrhoea, vomiting, dehydration)<sup>12</sup>. Amongst the major symptoms, watery diarrhoea characteristic of  
42 cholera is caused by the cholera toxin (CT)<sup>4-6</sup>. The *V. cholerae* El Tor responsible for the current cholera  
43 pandemic has become more virulent by undergoing several changes in CTX genotype<sup>13</sup> and acquiring  
44 virulence-related gene islands<sup>14</sup>.

45 In this study, we developed a reference-agnostic machine learning method, coupled with genome-scale  
46 metabolic modelling (GSMM) and protein structural analysis, to achieve two key objectives as outlined  
47 below. The first objective was to identify the genetic variations and signatures of the BD-1.2 lineage  
48 evolution beyond what has been found so far<sup>10</sup>. Our analysis considered 129 *V. cholerae* isolates from  
49 diarrhoea samples collected between 2015 and 2021, from patients admitted to the icddr,b hospital in  
50 Bangladesh. Several genomic studies investigated the evolution of lineages from 1991 to 2017, as well  
51 as in 2022<sup>9-11</sup>. However, there remains a gap in research during the intervening period. In our analysis,  
52 we discovered a set of 77 SNPs within the coding genome (mapped to 50 known genes), along with 12

53 annotated accessory genes, including some associated with antibiotic resistance, virulence, motility,  
54 colonization, biofilm formation, acid tolerance and bacteriophage resistance, identified as correlated  
55 with BD-1.2 transmission. Our findings go beyond what was recently discovered<sup>9-11</sup> for the lineage.

56 The second objective was to investigate if correlations exist between the genomic determinants of BD-  
57 1.2 strains and clinical manifestations among hospitalised patients from whom the isolates were  
58 collected from. Machine learning revealed the existence of correlations between genetic determinants  
59 in *V. cholerae* and clinical symptoms (diarrhoeal duration, number of stools, abdominal pain, vomit,  
60 and dehydration). Overall, the analysis revealed an overlap of 11 mutations, four accessory genes, and  
61 one intergenic SNP between the unique genomic determinants associated with BD-1.2 transmission and  
62 the clinical symptoms linked to this lineage. Additionally, a distinct set of 17 mutations, 39 accessory  
63 genes, and four intergenic SNPs were found exclusively linked to the severity of clinical symptoms.  
64 Through detailed GSMMs and 3D structure analysis of these genes, we inferred the mechanistic basis  
65 behind the selection of these genomic drivers in BD-1.2 and link to severity of the symptoms.

66

## 67 **Results**

### 68 **From 2015 to 2021 in Bangladesh, a diverse array of genetic variations characterises the** 69 **emergence of distinct circulating lineages**

70 To explore the evolutionary dynamics of *V. cholerae* linked to cholera cases in Bangladesh, a genomic  
71 analysis was done considering the years 2015 to 2021. We sequenced 129 *V. cholerae* O1 El Tor isolates  
72 taken from stool samples of patients between September 2015 to April 2021 admitted to hospitals in six  
73 districts (Barisal, Chittagong, Dhaka, Khulna, Rajshahi and Sylhet) of Bangladesh, Supplementary Data  
74 1. During the duration of this study, isolates belonging to serotypes Inaba and Ogawa were identified,  
75 Fig. 1. Consistent with previous studies<sup>10,15</sup>, a serotype switch was observed, with Inaba predominantly  
76 present in 2016 and 2017, followed by a predominance of Ogawa samples in 2018 and 2019 (Fig. S1).  
77 Both serotypes were detected in 2015 and continued to coexist from 2020 onwards. Serotypes were

78 significantly associated with collection years (chi-square test with p-value Bonferroni < 0.005) but not  
79 significantly associated with collection location (chi-square test with p-value Bonferroni > 0.005).

80 The maximum likelihood phylogeny of the 129 isolates was reconstructed based on the alignment of  
81 the core genome (3468 genes) and showed two distinctly evolved lineages, Fig. 1. Comparison with  
82 previous studies<sup>9,10</sup>, identified these lineages as BD-1.2 (n=84) and BD-2 (n=45), Fig. S2. Apart from  
83 the previously reported genetic variations<sup>4</sup>, we identified additional differences existing between the  
84 two lineages, in VSP (vibrio seventh pandemic; VSP-1 and VSP-2), VPI (vibrio pathogenicity islands,  
85 VPI-1 and VPI-2) and PLE (phage inducible chromosomal island-like elements), see Fig. 1. More  
86 precisely, in VSP-2, BD-2 isolates had a tryptophan at position 249, while BD-1.2 had a leucine at this  
87 position. In addition, in VSP-2, gene VC-514 (*aer*) was present in all BD-2 isolates but absent in BD-  
88 1.2. In VPI-2 a SNP led to an amino-acid variation at position 150, with BD-1.2 having an aspartic acid,  
89 and BD-2 an asparagine. BD-2 samples exclusively exhibited PLE2, while BD-1.2 samples had both  
90 PLE1 and PLE2 along with PLE2. Moreover, further differences were found in nonsynonymous SNPs  
91 on core genes and presence/absence of accessory genes, as described in the following section.

92 The distinct phylogeny patterns of BD-2 and BD-1.2, were also confirmed through a comparative study  
93 analysing 1134 isolates from *V. cholerae* El Tor O1 strains across 84 countries, including our isolates,  
94 (Supplementary Data 2 and 3, Fig. S3). BD-2 isolates clustered with Indian-1 (IND-1), while BD-1,  
95 BD-1.1, and BD-1.2 isolates from Bangladesh clustered with African (T9-T13)<sup>16</sup>, Latin America-3  
96 (LAT-3)<sup>13</sup>, Asian-2 (AS-2), and Indian-2 (IND-2) lineages (Fig. S3), in agreement with previous  
97 results<sup>10</sup>.

98

### 99 **Genetic and temporal differentiation of *V. cholerae* BD-1.2 and BD-2 lineages correlate with SNPs** 100 **on coding and non-coding regions, and accessory genes**

101 To assess the relatedness of *V. cholerae* isolates in our cohort, we measured the number of different  
102 core genome SNPs in a pairwise manner across all isolates. We created a network based on clusters of  
103 related isolates with less than 15 SNPs, as done previously<sup>17,18</sup>. Across the cohort the median SNP

104 difference was 117 SNPs (ranging from 0 to 1710 SNPs with IQR of 1211). The resulting undirected  
105 graph (Fig. 2) revealed that BD-2 and BD-1.2 formed two disconnected graphs each composed of  
106 samples from a specific lineage, but with no distinct separations between the Ogawa and Inaba  
107 serotypes.

108 To identify additional potential involvement of genetic elements in shaping the differences between the  
109 BD-1.2 and BD-2 isolates in our cohort, beyond current annotations (*ctxB* allele, type of SXT/ICE,  
110 VSP-II, VIP-I, *gyrA* gene allele)<sup>10</sup>, we sought for patterns of similarities and differences, at a finer scale,  
111 searching for the number, type and position of accessory genes as well as mutations in the core genome  
112 and intergenic regions across all the isolates. A two-sided Fisher exact test, with Bonferroni correction,  
113 was performed to assess the relationship between the BD-2 and BD-1.2 lineages and each of the various  
114 genomic features (core and intergenic SNPs and accessory genes). Overall, we found a significantly  
115 larger proportion of core genome mutations (51.4%, 1224 core genome SNPs and 73.1%, 160 intergenic  
116 SNPs) and a small proportion of accessory genes (11.3%, 115 genes) that exhibited statistically  
117 significant differentiation between the two lineages, Supplementary Data 4. Refer to Supplementary  
118 Note 1 and Fig. S4 for more details on the statistical analysis comparing the number of accessory genes,  
119 core genome SNPs and intergenic SNPs. The comparative analysis also indicated a temporal shift in  
120 the distribution of core genome and intergenic SNPs over the years, showing that BD-1.2 isolates  
121 accumulated different SNPs compared to BD-2 isolates as time progressed (Fig. S4E-F).

122 Out of the 115 accessory genes that differed between the two lineages, 12 were annotated while the  
123 remaining 101 were hypothetical. Among these 12 annotated genes, five – (*lon\_3*, *endA*, *adh*, *hdfR\_4*  
124 and *bcr\_2*) – were predominant (over 96% presence) in BD-1.2 and absent in BD-2, and seven (*aer\_3*,  
125 *hlyA\_2*, *mcrC*, *mepM\_3*, *mrr*, *tetA* and *tetR*) were present (over 97% presence) in BD-2 and absent in  
126 BD-1.2. Of the twelve annotated genes, three are known to be antimicrobial resistance genes (*bcr*, *tetA*  
127 and *tetR*)<sup>19</sup>. *TetA* and *tetR* were mainly detected in BD-2 isolates (97.7%), confirmed as primarily  
128 tetracycline-resistant through susceptibility testing in both doxycycline and tetracycline antibiotics  
129 (Supplementary Data 1). On the contrary, *bcr*, a multidrug efflux pump, was predominantly present in  
130 BD-1.2 isolates (96.4% of isolates) and completely absent in BD-2 isolates. Out of the 16 known

131 antimicrobial resistant genes (ARGs) present in the pangenome of this cohort, only *tetA*, *tetR* and *bcr*  
132 were found to statistically separate both lineages. *TetA* and *tetR* were both located in a contig showing  
133 high similarity to the SXT-ICE element, SXT(HN1) in BD-2 isolates. Conversely, *bcr* was found in a  
134 mobile element in the BD-1.2 isolates with similarity to SXT ICE element, ICEVchBan5. The presence  
135 of these SXT elements in the BD-2 and BD-1.2 lineages was previously shown by Monir *et al*<sup>10</sup>. Both  
136 contigs contained two identical insertion sequences, mobile genetic elements MGEs, (*ISShfr9* and  
137 *ISVsa3*), see Fig. S5. Also, among the 12 annotated genes, four (*endA*, *hlyA*, *lon* and *mcrC*) were  
138 previously found to be related to virulence<sup>18-23</sup>. More information about the function of these genes is  
139 given in the Supplementary Note 2.

140 To assess the extent of our results beyond our cohort, we investigated whether the 12 annotated  
141 accessory genes that we had found were also present in other Bangladeshi and Indian lineages. We  
142 performed a comparative genomic analysis of 219 *V. cholerae* O1 reference isolates collected in  
143 Kolkata, India, and Dhaka, Bangladesh, between the years 2004 and 2022 (ENA public database  
144 <http://www.ebi.ac.uk/ena>, see Supplementary Data 5). The results confirmed the presence/absence  
145 patterns of the 12 genes in the BD-1.2 and BD-2 lineages in the reference isolates, aligning with our  
146 initial findings, see Supplementary Note 2.

147 In addition to differences in accessory gene types and patterns, missense mutations associated to allelic  
148 variations were found in BD-1.2, when compared to BD-2 strains. We identified 1385 SNPs in the core  
149 genome, including 291 non-synonymous and 934 synonymous coding variants, both representing  
150 variants in their functional protein-coding form. In addition, 160 intergenic SNPs were found,  
151 representing variants in their regulatory form. Many SNPs showcased unique allelic distribution  
152 patterns between the two lineages. When mapped back, the non-synonymous SNPs identified 291  
153 amino acid substitutions in 105 genes, including 50 known genes and 55 hypothetical ones (see  
154 Supplementary Data 4). Table S1 shows core genes with allelic distribution between BD-1.2 and BD-2  
155 significantly different (i.e., containing polymorphic sites found exclusively in one lineage but absent in  
156 the other lineage).

157

158

159 Among the genes exhibiting lineage-specific allelic variation, some contribute to functions including  
160 growth, cell wall organization, colonization, toxigenicity and resistance, similar to what found  
161 previously<sup>10</sup>. Additionally, we found genes with a unique non-synonymous variant in BD-1.2, with roles  
162 in toxin transport and acid tolerance, shedding light on functions that may clarify their contribution to  
163 the recent prevalence of BD-1.2 over BD-2. See Supplementary Note 3 for more information about  
164 these genes. Notably, *ompU* is another gene with a statistically significant mutation (G325D) underlying  
165 lineages' separation. Amino acid D is predominant in BD-1.2, while the amino-acid G is prevalent in  
166 BD-2. To assess for any additional genes separating the BD-1.2 and BD-2 lineages we also conducted  
167 an analysis on the pangenomes of the lineages separately but found the results broadly in line with that  
168 of the combined pangenome analysis presented above (Supplementary Note 4 and Supplementary Data  
169 6-10)

170 To understand the systemic relationships connecting the identified lineage-specific genetic signatures  
171 on a mechanistic level, we analysed the 30 core genes in Table S1 with allelic variants that were found  
172 exclusively in one lineage but absent in the other lineage using the *V. cholerae* GSM model iAM-Vc960  
173 (Fig. 4). Thirteen of these genes (*murI*, *ftsI*, *appC*, *suhB*, *glmM*, *dsbD*, *licH*, *cysG\_1*, *cobB*, *clcA*, *argG*,  
174 *mak*, *phhA*) are metabolic and have been identified as playing integral roles in amino acid metabolism,  
175 cell wall metabolism, carbon metabolism, amino sugar and nucleotide sugar metabolism, energy  
176 metabolism (see Supplementary Data 11). Moreover, for these genes we sought to better understand  
177 their role by examining their effects on *V. cholerae* growth rate biochemical networks and metabolites  
178 production in the networks. As the effect of mutations/gene knockouts cannot always be observed as  
179 change in growth rate (due to the redundancy of the reactions in metabolic networks of bacteria), it can  
180 be useful to also consider the changes in metabolite yield. Changes in metabolite yield have been found  
181 to correlate with changes in the virulence, persistence, and fitness of some organisms<sup>24</sup>. Furthermore,  
182 *V. cholerae* are capable adapting to ecological niches by altering the metabolites they excrete to create  
183 a more favourable environment for *V. cholerae* and/or a less favourable environment for other species  
184 competing for the same resources<sup>25,26</sup>. Mutations disrupting larger numbers of metabolite yields may be



185 suggestive of a larger systems-level impact on bacterial metabolic function. Therefore, gene  
186 essentiality, flux variability analysis (FVA) and flux balance analysis (FBA) were used to predict,  
187 through gene knockouts, the essentiality and the effects of the identified genetic determinants on the  
188 growth rates of *V. cholerae*, and also used to further explore their influence on metabolite yield. The  
189 latter was done by assessing the influence on metabolite flow within the complete metabolic network  
190 of *V. cholerae*, encompassing all known metabolites and metabolic reactions (see Methods). In this  
191 analysis it was important to consider all reactions and metabolites in the model rather than focussing on  
192 a subset, as doing so ensures no undue bias or assumptions underlie the results.

193 The genes *cysG*, *clcA*, *adh* and *mcrC*, were found to be essential for growth (i.e., knocking these genes  
194 out reduced the biomass growth to less than  $0.0001\text{h}^{-1}$ ) in both rich and minimal media. Furthermore,  
195 *murI*, *glmM*, and *dapF* displayed auxotrophic behaviour in minimal media, whereas *cysG*, *clcA*, *adh*,  
196 and *mcrC* were found to be essential in rich media with alternative carbon sources. Additionally, three  
197 genes, *murI*, *glmM* and *dapF*, were found to be essential for growth in minimal media only. Next, flux  
198 variability analysis (FVA) was used to identify biochemical reactions whose flux span was significantly  
199 changed (greater than 10% change) by knocking out these genes. In total ten genes *murI*, *glmM*, *cysG*,  
200 *clcA*, *argG*, *mak*, *adh*, *dapF*, *add*, and *mcrC* when knocked out significantly changed the flux span in  
201 at least one reaction through the model by FVA analysis, Supplementary Data 11. Finally, FBA analysis  
202 was used to determine the effect of gene knockouts on metabolite yield. Five genes, *murI*, *glmM*, *cysG*,  
203 *mak*, and *dapF* were found to reduce at least one metabolite yield to zero in the model when knocked  
204 out (given the wildtype yield was greater than 0), Supplementary Data 11<sup>27,28</sup>. Interestingly, the average  
205 number of metabolite yields affected by knockouts of the genes discriminating lineages was  
206 significantly higher than a random selection of 100 metabolic genes (p-value 0.0429, Mann Whitney U  
207 test, two-sided), indicating a stronger influence on metabolite production for this subset of genes.

208 To further elucidate the metabolic differences between the BD-1.2 and BD-2 lineages, we repeated our  
209 previous analyses done on the generalized model using strain-specific models automatically generated  
210 by CarveMe<sup>27</sup>. Gene essentiality analysis concurred with the general model (iAM-Vc960), with only a  
211 small number of differences (Supplementary Data 12). The effect of *murI* gene knockouts differed  
212 between lineages, proving non-essential in 94% of BD-1.2 lineage models but only in 76% of BD-2

213 lineage models. Flux variability analysis of the individual models revealed that *clcA* knockouts led to  
214 significant changes in the flux span of the CLt3\_2pp reaction, which controls chloride transport, in 96%  
215 BD-2 models compared to just 5% of BD-1.2 models. The *clcA* gene has been linked to bacterial acid  
216 resistance and it has been suggested that changes to the expression/repression of this gene may help  
217 facilitate survival during movement through the intestinal tract <sup>28</sup>. Similarly, flux balance analysis  
218 indicated that metabolite yield was changed differently across lineages in response to knocking out  
219 *clcA*, with the metabolite yield of chloride reduced to 0 in 95% of BD-1.2 isolates.

220 In summary, a total of 15 genes found to underly the genetic and temporal differentiation of *V. cholerae*  
221 BD-1.2 and BD-2 lineages, were also found to significantly alter the growth, reaction flux, or metabolite  
222 yield of *V. cholerae* when knocked down, either in the generalised iAM-Vc960 GSM model or in the  
223 draft strain-specific models. Of interest was the gene *clcA*, which showed differences in both flux span  
224 and metabolite changes between lineages in the draft GSM models. The FVA and FBA results indicate  
225 that the genes identified by machine learning as strongly associated with the severity of symptoms play  
226 important metabolic roles. Disruption of these functions could potentially affect bacterial growth or  
227 metabolic output, which may contribute to the survival and dominance of one lineage over another.  
228 Although our analysis cannot pinpoint a single SNP as responsible for the loss of metabolic function, it  
229 suggests that an accumulation of SNPs or gene losses could collectively lead to metabolic changes. We  
230 observe the potential for metabolic alterations driven by multiple mutations (SNPs).

231 Lastly, when mapping the 160 intergenic SNPs back to genomes, we found their location in the  
232 upstream/downstream regions of 35 known genes and 34 hypotheticals genes (see Supplementary Data  
233 4). These intergenic SNPs exhibited allelic distribution, with the minor variant prevalent in the BD-2  
234 isolates (68% to 100%), while the major variant dominated in the BD-1.2 isolates (over 98%), only one  
235 SNP in BD-1.2 had a major allelic variant at of 47% (Fisher exact test, Bonferroni correction p-value<  
236 2.31e-08). Many of these SNPs were located within transcriptional factor binding sites (TFBs)  
237 (Supplementary Data 4). Intergenic SNPs, exhibiting significantly different allelic distributions  
238 between BD-1.2 and BD-2, mapped across the TFBs of 11 TFs (*ToxT*, *Fur*, *AmpR*, *OmpR*, *LuxR*, *LexA*,

239 *ArgR*, *PhoP*, *CRP*, *ArcA*) (Fig. S6-S16). More information about the function of these transcriptional  
240 factor binding motifs is provided in Supplementary Note 6.

241

## 242 **Machine learning unravels correlations between genomic determinants and clinical symptoms in** 243 **humans**

244 Beyond identifying the potential involvement of new genetic traits in differentiating the BD-1.2 and  
245 BD-2 lineages, we hypothesized that the same or additional genetic features might play a significant  
246 role in the manifestation and severity of clinical symptoms in patients when infected with *V. cholerae*.

247 A summary of the distribution of each clinical symptom over the two lineages is given in Fig. S17. We  
248 focused on the lineage BD-1.2, which caused the most recent outbreak in Bangladesh. To identify if  
249 and which coding and non-coding mutations and/or presence/absence of accessory genes would  
250 correlate with the different clinical symptoms, we employed a bespoke, supervised machine learning  
251 pipeline.

252 The pipeline is aimed at mining sequencing data to identify the genetic elements that more strongly  
253 correlate with observed clinical symptoms, which in this case are vomit, dehydration, number of stools,  
254 duration of diarrhoea and abdominal pain (see Methods section). The pipeline is a bespoke adaptation  
255 of ML-based data-mining methods previously developed within our team to identify correlations  
256 between genomic features with phenotypes<sup>17,18,29,30</sup>. In the pipeline, information about different genetic  
257 features (SNPs -both from coding and non-coding regions- and presence/absence of accessory genes)  
258 can be encoded as input to ML-powered predictive models designed to estimate the likelihood of  
259 observing the selected phenotypes under each specific pattern of input values<sup>17</sup>. As long as trained with  
260 sufficient observational data, the ML-powered predictive models are able to replicate experimental  
261 evidence, in addition to providing information on what inputs correlated most strongly with each  
262 phenotypic manifestation. Through such introspective power, the pipeline is able to unravel co-  
263 occurrent, multiple mechanisms (mutations, horizontal gene transfer - HGT), variants in their functional

264 protein-coding and regulatory forms, as well as their additive effect on the targeted phenotypes, which  
265 in this work, were clinical symptoms.

266 The following clinical symptoms were selected, namely: vomit, abdominal pain, diarrhoea duration,  
267 24-hour stool count and dehydration. Each clinical symptom was handled by building a dedicated  
268 symptom prediction model, operating using genetic elements as inputs. Two symptoms (vomit and  
269 abdominal pain) were encoded as binary (presence vs absence). The other three symptoms – diarrhoea  
270 duration, 24-hour stool count, and dehydration – were encoded as multi-class: dehydration as None,  
271 Moderate and Severe; diarrhoea duration as < 1 day, 1-3 days, 4-6 days, and 7-9 days; and stool count  
272 in 24 hours as 3-5 times, 6-10 times, 11-15 times, 16-20 times, and 21+ times. We handled the prediction  
273 of multi-class symptoms via the implementation of binary predictors.

274 The symptom prediction models were developed with built-in robustness to potential confounding  
275 factors. Specifically, the following list of variables was initially considered as potentially having  
276 confounding effects: year of collection, location of patient, sex of patient, age of patient and serology  
277 of *V. cholerae*. Each potential confounder was tested for correlation to the symptom being targeted by  
278 the prediction model. If the potential confounder was found correlated to the symptoms (hence moving  
279 from potential to proven confounder), then any other input variable also found correlated with the same  
280 confounder would be eliminated from the prediction model. All the correlation tests between inputs  
281 and symptoms, as well as between inputs themselves, were run using two-sided Chi-square tests.  
282 Further, possible confounding effects related to random initialisation parameters of SMOTE (see  
283 methods) were contained by running SMOTE multiple times.

284 The development and optimisation of each symptom prediction model powered by machine learning  
285 was based on running a comparative analysis of the predictive performances of different machine  
286 learning algorithms, namely: linear support vector machine (linear SVM), non-linear SVM with radial  
287 basis function (RBF SVM), random forest, extra-tree classifier and logistic regression) and two meta-  
288 methods (Adaboost and XGBoost). For each algorithm, multiple configurations of the hyperparameters  
289 of the learning algorithms were tested. A nested cross validation approach was used to select the best  
290 hyperparameters, based on randomly selecting different training and test sets, and using stratified k-

291 fold cross validation metric. Finally, Friedman and Nemenyi tests were used to statistically compare  
292 and select the best performing algorithm for each prediction model (see Methods section).

293 In the end, based on a two-sided Chi-square test of independence (p-value < 0.01), the models for  
294 abdominal pain, vomit, number of stools 11-15 times vs. 21+ times, number of stools 11-15 times vs.  
295 16-20 times, dehydration moderate vs severe were found immune to confounding effects due to year of  
296 collection, location of patient, sex of patient, age of patient and serology of *V. cholerae*. The prediction  
297 model: diarrhoea duration <1day vs 1-3 days was found immune to confounding effects due to age of  
298 patient, sex of patient, location of patient, and serology of *V. cholerae*. However, the prediction model  
299 was found to be influenced by year of collection; therefore, the inputs that were also correlated to year  
300 of collection were removed from the analysis (Supplementary Data 13). Moreover, we were able to  
301 successfully develop six binary symptom prediction models featuring adequate prediction performance  
302 levels. These were dedicated to predicting the following binary phenotypical outcomes: i) stools 11-15  
303 times vs. 16-20 times; ii) stools 11-15 times vs. 21+ times; iii) moderate vs. severe dehydration; iv)  
304 diarrhoea duration <1 day vs. 1-3 days; v) presence vs absence of vomit; and vi) presence vs absence  
305 of abdominal pain (Supplementary Data 14). The remaining binary predictors were discarded for not  
306 performing adequately, either because of unbalanced available sets of observations (needed for training  
307 the supervised ML models), or because of more challenging separability of the phenotypes given the  
308 selected inputs (no features were statistically significant based on the Fisher exact test). Among the  
309 tested pipeline technologies mentioned earlier, logistic regression was identified by the Friedman F-test  
310 and the Nemenyi post-hoc analysis as the best performing one (Fig. S18). Of the six binary prediction  
311 models, four had an AUC greater than 0.9, Fig. 4. Supplementary Data 15 indicates the performance  
312 metrics obtained by all binary predictors for each clinical symptom. Figs. 4 and S19 show the  
313 performance results for the Logistic regression classifier.

314 Analysis of the best-performing symptom prediction models allowed us to identify the input features  
315 (core genome coding and intergenic SNPs and accessory genes) most strongly correlated to each  
316 phenotype (Supplementary Data 16). Seventy-nine different features in total were selected as  
317 significantly correlated to at least one of the six symptom prediction models, with 68% being selected

318 in two or more models (Fig. 5). No features were selected for all symptoms. All features associated with  
319 number of stools 11-15 times vs. 21+ times were found associated to at least one of the other five  
320 symptom prediction models. Forty-five accessory genes (nine known genes, *tufB\_2*, *blc*, *pckA*, *luxR\_2*,  
321 *hcpA\_1*, *rpoS*, *dcuA*, *hpt*, *luxR*, and 36 hypothetical genes) and 28 core SNPs over 23 genes (14 known,  
322 *clpS*, *gshB*, *dapF*, *fabV\_1*, *add*, *tufB*, *lpoA*, *phrB*, *yjcS*, *fabH1*, *cysG\_2*, *padC*, *pepN*, *tadA\_2*, and nine  
323 hypothetical genes) were identified as strongly associated to at least one of the symptoms. From the  
324 nine known accessory genes: four (*rpoS*, *hpt*, *luxR* and *pckA*) were found in the vomit model; *dcuA* was  
325 found in the abdominal pain model; *hcpA\_1* was found only in the number of stools 11-15 times vs. 16-  
326 20 times; *luxR\_2* was found in two models (vomit and dehydration moderate vs severe); *blc* and *tufB\_2*  
327 were found in three models (vomit, number of stools 11-15 times vs. 16-20 times and number of stools  
328 11-15 times vs. 21+ times) with *tufB\_2* also found in abdominal pain and diarrhoea duration <1 day  
329 vs. 1-3 days models. Six SNPs from the genes *tufB*, *dapF*, *clpS*, *gshB* and *fabV* were associated to three  
330 symptom prediction models (vomit, number of stools 11-15 times vs. 16-20 times and number of stools  
331 11-15 times vs. 21+ times) with the SNPs from the genes *dapF* and *fabV* also associated with abdominal  
332 pain and diarrhoea duration <1 day vs. 1-3 days and the SNP from the gene *tufB* associated with  
333 dehydration moderate vs severe.

334 Among the 45 accessory genes linked to clinical symptoms, six hypothetical genes were also  
335 statistically significant in distinguishing the two lineages. Among the other accessory genes selected,  
336 four (*blc*, *pckA*, *luxR* and *rpoS*) have important biological functions. In particular, *Blc*, also known as  
337 *VlpA*, is a lipocalin, that is correlated to acquisition of drug resistance in *V. cholerae*<sup>31</sup>. *PckA*  
338 (phosphoenolpyruvate carboxykinase) is important for gluconeogenesis, a highly conserved pathway in  
339 bacteria and humans. Interfering with gluconeogenesis pathway impacts *V. cholerae* colonization in  
340 mouse models, highlighting its crucial role in sustaining *V. cholerae* growth and viability within the  
341 intestines<sup>32</sup>. *LuxR* plays a key role in regulating biofilm production and secretion in *V. cholerae*<sup>33</sup>. *RpoS*  
342 is a sigma factor that facilitates physiological adaptation to general starvation and stationary phase  
343 growth in different species. *V. cholerae* strains lacking the gene *rpoS* are impaired in the ability to

344 survive in different environmental stresses. *RpoS* was also shown to be important in *V. cholerae* for  
345 efficient intestinal colonization<sup>34</sup>.

346 Out of the 28 core SNPs associated to the clinical symptoms, 11 were also found previously as  
347 statistically significant in differentiating the BD-2 and BD-1.2 lineages (see above), Supplementary  
348 Data 16. These 11 SNPs mapped to 11 genes (*clpS*, *gshB*, *dapF*, *fabV\_1*, *add*, and six hypothetical).  
349 Among the SNPs mapping to known genes (*clpS*, *gshB*, *dapF*, *fabV\_1*, *add*), three are non-synonymous  
350 SNPs mapping to *clpS*, *gshB* and *fabV*. In *V. cholerae* ClpS regulation involves cAMP receptor protein  
351 (CRP)<sup>31</sup>. CRP is important in intestinal colonization<sup>35</sup>. *GshB*, encodes a glutathione synthetase (GSH),  
352 a gene associated to resistance to oxidative stress. *V. cholerae fabV* is one of the several triclosan-  
353 resistant ENR encoding genes<sup>36</sup>.

354 As in our previous lineage analysis, we sought to better understand the importance of the genes which  
355 had been found to better correlate with the severity of the symptoms. We examined for those genes that  
356 were metabolic, through FVA and FBA, the effects of such genes on growth rate (gene essentiality),  
357 and beyond that, their influence on metabolite yield and reaction flux. Nine symptoms-related genes  
358 were identified as metabolic genes in the iAM-Vc960GSM model (Fig. 6). Eight of these genes were  
359 associated to five metabolic systems Supplementary Data 17). *FabH1* and *gshB* associated with  
360 cofactor and prosthetic group metabolism; *pckA* is associated with carbohydrate metabolism; *dcuA*  
361 plays a crucial role in C4-dicarboxylate transport; *dapF*, *pepN* and *gshB* are significant in amino acid  
362 metabolism; *add* and *pckA* are relevant to nucleotide metabolism; *oppA* and *fabH1* are involved in cell  
363 wall metabolism, with *fabH1* relevant for fatty acid biosynthesis (Supplementary Data 17).

364 Using FBA and FVA analysis, the knockouts of the genes *dapF* and *gshB* were found to halt production  
365 of several metabolites. The genes *pckA*, *add*, *dapF*, *oppA*, *gshB* were found to significantly change the  
366 reaction flux span, Supplementary Data 17. Both FBA and FVA analysis can infer if potential metabolic  
367 adaptation mechanisms for *V. cholerae* can lead to alterations in bacterial virulence, potentially leading  
368 to worst symptoms, if genes significantly affect pathways which are associated to important functions  
369 such as colonization, biofilm production and cell wall synthesis. For example, the *gshB* gene, a

370 glutathione reductase, contributes to *V. cholerae* intestinal colonization<sup>37</sup> and has a role in acid tolerance  
371 response<sup>38</sup>. Similarly, *dapF* was found as an essential gene in minimal media and leading to auxotrophic  
372 behaviour to the amino-acid lysine. As Pearcy *et al.*<sup>39</sup> indicated, an auxotrophic behaviour of a gene  
373 connected to amino-acid biosynthesis is important because it can provide competitive fitness advantage  
374 against commensal bacteria. During the infection stage *V. cholerae* engage and compete with  
375 commensal bacteria for nutrient acquisition to support rapid growth and multiplication<sup>40</sup>. Moreover, the  
376 lysine pathway plays a central role in eubacteria cell wall biosynthesis, since meso-diaminopimelate is  
377 the immediate precursor for the biosynthesis of its main component, peptidoglycan, with *dapF*  
378 responsible for the creation of meso-diaminopimelate in the lysine pathway<sup>41,42</sup>. The proper synthesis  
379 and maintenance of peptidoglycan is essential for bacterial virulence and its viability<sup>43</sup>.

380 To further investigate the link between metabolic gene variations and the clinical symptoms observed  
381 in different strains, we utilized draft strain-specific models generated with CarveMe<sup>27</sup>. The gene  
382 essentiality analysis results were largely consistent with those of the general model (iAM-Vc960), with  
383 only a few differences noted (Supplementary Data 18). The effect of *dapF* gene knockouts varied  
384 between models with the gene being essential in 93% (n=20) and non-essential in 7% (n=9) of the  
385 models. Comparing symptoms between the ‘essential’ and ‘non-essential’ groups, dehydration was  
386 significantly more severe in the ‘non-essential’ group (Fisher exact test *p* value =0.05). All strains in  
387 this group exhibited severe dehydration, suggesting a link between non-essentiality of the *dapF* gene  
388 and the severity of *V. cholerae* symptoms. In relation to this, the flux balance analysis revealed changes  
389 in metabolite yields associated with the genes *dapF* and *cysG\_2* across all strain-specific models. For  
390 *dapF*, altered metabolite yields were predominantly observed in strains where *dapF* was essential, while  
391 knocking out *dapF* in non-essential models had minimal impact on the metabolite yields of murein-  
392 related metabolites. This indicates metabolic adaptations linked to bacterial survival in these strains,  
393 potentially contributing to more severe disease outcomes. Additionally, knocking out the *padC* gene  
394 resulted in significant changes in metabolite yields only in the NGICDV-066 strain. Although  
395 conclusions drawn from a single strain are limited, it is notable that this isolate exhibited the most severe  
396 clinical symptoms across all measured symptoms, except for the duration of diarrhoea (presence of



397 vomiting, presence of abdominal pain, number of stools (21+ times), presence of severe dehydration,  
398 duration of diarrhoea 1-3 days). Flux variability analysis in individual models indicated consistent  
399 behaviour across all strain-specific models regarding gene knockouts associated with clinical  
400 symptoms. Specifically, five gene knockouts (*add*, *dapF*, *gshB*, *padC*, *pckA*) showed significant flux  
401 span changes in all models.

402 In summary, in relation to gene essentiality, reaction flux and metabolite yield, our results show that  
403 *gshB* and *dapF* make interesting candidates for further analysis, as knockout models of these genes  
404 predict significant changes to the bacterial metabolic function.

405 To delve deeper into understanding the functional mechanisms underlying clinical symptoms, we  
406 explored the interactome of the proteins associated to the clinical symptoms. The protein-protein  
407 interaction network (PPI) analysis revealed the interactome of 36 proteins, selected by the machine  
408 learning pipeline, with 109 other proteins, Fig. S20. The KEGG analysis indicated enrichment in  
409 ribosome proteins (e.g., RpoS) and fatty acid biosynthesis (e.g., FabH1, FabV) (Fig. S21). The  
410 colonization in the human intestine and virulence of *V. cholerae* is intricately connected to both fatty  
411 acid metabolism<sup>44</sup> and the ribosome pathway<sup>45</sup>. The GO analysis highlighted enrichment in translation,  
412 peptide biosynthetic processes, and gene expression, featuring TufA, TufB, RpoS, GshB  
413 (Supplementary Data 19 and 20). The peptide biosynthetic pathway plays a vital role in *V. cholerae*  
414 biofilm formation and colonization<sup>23</sup>.

415 None of the six intergenic SNPs selected by the machine learning pipeline were in TFBs or promoters.  
416 These SNPs were located in a region without any functional annotations within 2 kbps upstream or 0.5  
417 kbps downstream of a gene, adhering to the standard database dbSNP cutoffs for SNP-to-gene  
418 mapping<sup>46,47</sup>. See Supplementary Data 16 for additional information about the location of these SNPs.

419

420 **Structural analysis suggests evolutionary drivers of selection, mechanistic bases for BD-2 and BD-**  
421 **1.2 lineages evolution, and associations to clinical symptoms**

422 To further understand whether the identified alleles play a causal role in the evolution of lineages and  
423 clinical symptoms, we selected two of the top-ranked non-synonymous SNP candidates, prioritizing the  
424 following aspects in relation to the associated genes: (i) have significant difference of allelic distribution  
425 between BD1-1.2 and BD-2; (ii) have a significant correlation, as detected by the ML pipeline, with the  
426 selected clinical symptoms; (iii) are characterised as functionally important for *V. cholerae* metabolisms  
427 (i.e. significantly impacting reaction flux when knocked out, as highlighted by the GSM model) and/or  
428 interactome (i.e. enrichment of the functions and mechanisms related to pathogenesis); (iv) 3D  
429 structural mutation analysis could be benchmarked with experimental evidence. This resulted in three  
430 genes, all top-ranked by both the Fisher Exact test for BD-1.2 and BD-2 lineage evolution and the ML  
431 analysis for the underlying clinical symptoms, namely: *fabV*, *gshB* and *clpS*. We mapped the alleles of  
432 *fabV*, *gshB* and *clpS* to their protein structures using both experimental crystal structures and predicted  
433 homology models. However, the 3D-structure could be utilised to infer the mechanistic basis only for  
434 *fabV* and *gshB*.

435 In all BD-2 isolates FabV had a proline at position 149 (Pro149) whereas, in BD-1.2 isolates, the Pro149  
436 was found in only 40.5% of cases, with the remaining 59.5% isolates exhibiting histidine at position  
437 149 (His149). The BD-1.2 isolates with His149 showed a higher duration of diarrhoea (1-3 days) and a  
438 higher number of stool score (16-20 times and 21+ in 24 hours) compared to the BD-1.2 isolates with  
439 Pro149, featuring a lower diarrhoea duration (<1 day) and lower number of stools score (11-15 times).  
440 The amino acid 149 was located in the trans-2-enoyl-CoA reductase catalytic domain (Fig. 7A-E), when  
441 Pro149 is present, it interacts with Lys148, Ser151, Trp159 through Van der Waals (VDW) interactions,  
442 whereas His149 not only forms the aforementioned interactions but also creates an extra VDW  
443 interaction with Lys148. Furthermore, His149 interacts with an additional amino acid, Arg150, through  
444 a VDW interaction. These additional interactions in the presence of the His149 cause an increase in the  
445 stability of the structure ( $\Delta\Delta G = 0.101$  kcal/mol >0) and a decrease of the molecule flexibility ( $\Delta\Delta S_{\text{vib}}$   
446 ENCoM:  $-0.053$  kcal.mol<sup>-1</sup>K<sup>-1</sup>), which is usually linked to a stronger binding affinity<sup>48,49</sup>. Moreover, the

447 presence of His149 increased the positive charge of the surrounding area (Lys148, His149, Arg150)  
448 (Fig. S22), with an overall electrostatic energy increasing from 7.3E+03 kJ/mol (Pro149) to 7.48E+03  
449 kJ/mol (His149) within the 5Å region and with an overall protein total electrostatic energy rising from  
450 2.1E+05 kJ/mol (Pro149) to 2.52E+05 kJ/mol (His149). Exposed, positively charged amino acids are  
451 suggested to promote interactions with negatively charged cellular systems<sup>50</sup>. The enhanced positive  
452 charge of FabV in the presence of His 149 might support its role in participating in the breakdown of  
453 the negatively charged fatty acids.

454

455

456 GshB, a glutathione reductase, has been shown to contribute to *V. cholerae* intestinal colonization<sup>37</sup> and  
457 to have a role in the ability of *V. cholerae* to mount an acid tolerance response<sup>38</sup>. In all BD-2 isolates  
458 GshB had a threonine at position 93 (Thr93), whereas in the BD-1.2, the Thr93 was only found in 21.5%  
459 of the cases, with most (78.5%) of the BD-1.2 isolates exhibiting an isoleucine (Ile93) at this position.  
460 The BD-1.2 isolates with Ile93 are associated to a higher duration of diarrhoea (1-3 days) and a higher  
461 number of stool score (16-20 times and 21+ in 24 hours) compared to the BD-1.2 isolates with Thr93.  
462 Thr93 interacts with Asp92, Ile96, Tyr97 through 13 VDW interactions and 1 H-bond; whereas Ile93  
463 not only forms the aforementioned interactions but also creates extra VDW interactions with Tyr97  
464 (Fig. 8A-E). These additional bonds in the presence of Ile93 cause an increase in the stability of the  
465 structure ( $\Delta\Delta G = 0.384$  kcal/mol  $>0$ ) and a decrease of the molecule flexibility ( $\Delta\Delta SVib$  ENCoM: -  
466 0.055 kcal.mol<sup>-1</sup>.K<sup>-1</sup>), which is usually linked to a stronger binding affinity<sup>48,49</sup>. Moreover, the presence  
467 of Ile93 increased the negative charge of the surrounding area ( $<5\text{\AA}$ ) (Fig. S23A-B), with an overall  
468 electrostatic energy decreasing from 7.93E+03 kJ/mol (Thr93) to 7.4E+03 kJ/mol (Ile93) within the 5Å  
469 region and with an overall protein total electrostatic energy varying from 2.1E+05 kJ/mol (Thr93) to  
470 1.8E+05 kJ/mol (Ile93). A decrease in total electrostatic energy is often associated to folding<sup>51</sup>, protein  
471 folding stability is largely dependent on the hydrophobic interactions of nonpolar residues<sup>52</sup>. The  
472 surface, on average, has become more hydrophobic, indicating a possible reorientation of residues or a  
473 change in the surface's exposure to the solvent (Fig. S23C-D).

474

## 475 **Discussion**

476 Bangladesh has witnessed the continual genomic evolution of *V. cholerae* lineages, with increased  
477 virulence, resistance, global spreading ability and disease severity. The potential of a *V. cholerae* isolate  
478 to have a global spreading ability and cause disease is mostly approached by studying its genomics via  
479 bioinformatics analysis. Two recent studies<sup>9,10</sup> explored the genomics attributes of the lineage BD-2  
480 predominant between 2004 and 2018 and the emergent lineage BD-1.2 appearing from 2016 onwards  
481 and responsible for the 2022 outbreak<sup>9,10</sup>. By comparing these lineages, the authors revealed mutations  
482 in *ctxB* allele, SXT/ICE, VSP-II, VPI-1 and *gryA* allele<sup>10</sup> potentially explaining the recent shift in  
483 lineage predominance. Despite these knowledge advances, gaps persist in understanding the entire  
484 genomic repertoire associated to transmission ability and different disease severity patterns.

485 Here, we developed an analysis approach that combines, ML-powered data mining, whole-genome  
486 sequencing, genome-scale metabolic modelling and 3D structural analysis to uncover, on a finer scale,  
487 unknown associations between lineage transmission dynamics, diseases severity and the genomic make-  
488 up of *V. cholerae* isolates. Machine learning offers a powerful opportunity to analyse entire genomes  
489 efficiently against selected phenotypes (lineages, clinical symptoms), allowing for the identification of  
490 genomic features ranked on strength of correlation with the phenotype. This provides a significant  
491 advantage to conventional genomics-only methods based on checking for presence/absence or based on  
492 similarity searches of known manually chosen determinants. Moreover, our approach allowed various  
493 genetic determinants (accessory genes, and core coding and intergenic SNPs) to be analysed  
494 simultaneously to capture the co-occurrence, synergism and additive effect of multiple mechanisms and  
495 determinants (mutations, accessory genes, horizontal gene transfer, functional, metabolic, and  
496 regulatory variants). Determinants identified by ML may contain genes with a known functional  
497 relationship with the phenotype as well as genes with no previously known association with that specific  
498 phenotype. Altogether, our reference-agnostic approach overcomes limitations of previous genomics

499 studies that only considered one feature type (SNPs, accessory genes) at a time and known genetic  
500 elements associated to *Vibrio* transmission.

501 Using our method, in addition to confirming the aforementioned mutations identified in recent genomics  
502 studies<sup>10</sup>, we found further mutations in VSP, VPI, and PLE, exclusive to one lineage and absent in the  
503 other, supplementing those previously found by Monir et al.<sup>10</sup>. Moreover, our findings expand known  
504 mutations to a wider range of genomic determinants, including 115 accessory genes, 1225 core coding  
505 SNPs, and 160 intergenic SNPs crucial for explaining at a more-in depth scale BD-1.2 and BD-2 recent  
506 shift. Supplementing the previous knowledge on the type, number and functions of genomics  
507 determinants differentiating BD-1.2 and BD-2<sup>10</sup>.

508 For example, five core genes (*skp*, *tamA*, *clcA*, *cysG*, and *vals*) with a unique non-synonymous variant  
509 in BD-1.2 and playing key roles on toxin transport and acid tolerance, shed new light on functions and  
510 may help clarify their contribution to the recent prevalence of BD-1.2 over BD-2. In addition, non-  
511 synonymous SNPs, found uniquely in BD-1.2, were mapped to genes with functions such as  
512 colonization, toxins export, virulence, growth, response to pH and temperature, and phage resistance.  
513 For example, the mutation G325D in *ompU* conferring bacteriophage resistance<sup>29</sup>, was found in this  
514 work to be statistically important to differentiate the two lineages. OmpU a pore-forming protein of the  
515 outer membrane of *V. cholerae* has adhesive properties which may play a role in the pathogenesis of  
516 cholera<sup>53</sup>, is critical for vibrio fitness<sup>54,55</sup>, for dissemination<sup>54</sup>, for protection against the bactericidal  
517 effect of bile salts<sup>56</sup>, cationic peptides<sup>57</sup> and intestinal organic acids<sup>58</sup>. The G325D mutation is located  
518 within the L8 loop, which has been reported to be crucial for neutralizing infection and conferring  
519 resistance against phages<sup>59,60</sup>. Seed et al.<sup>60</sup>, showed that in presence of the bacteriophage ICP2  
520 (bacteriophage that preys on *V. cholerae* and was first isolated from cholera patient stool samples<sup>61</sup>) the  
521 OmpU virulent mutant (G325D) had a 10,000-fold enrichment over the wild-type, indicating that strong  
522 selective pressure is imposed by phage predation during *V. cholerae* infection.

523 Out of the twelve accessory genes found statistically significant to differentiate the two lineages, five  
524 (*lon\_3*, *endA*, *adh*, *hdfR\_4* and *bcr\_2*) were present uniquely in BD-1.2 with functions such as antibiotic  
525 resistance and biofilm formation. Increasing evidence indicates that *V. cholerae* has the capability to

526 develop biofilm-like aggregates during infection, potentially serving as a function in pathogenesis and  
527 disease transmission. Nonetheless, the composition, control mechanisms governing the formation of  
528 these biofilms during infection, and their significance in intestinal colonization and virulence remain  
529 yet to be elucidated<sup>62</sup>.

530 In addition to the coding genome, we found that regulatory networks are associated to lineage  
531 differentiation. Among the most relevant intergenic SNPs exhibiting significant allelic distribution  
532 between the two lineages is the one mapping in the TFBs of *ToxT*. This TF plays a crucial role in the  
533 development of *V. cholerae*-related symptoms<sup>60</sup> and selectively regulates the expression of virulence  
534 genes found in toxin-coregulated pilus (TCP) and cholerae toxin (CT)<sup>63,64</sup>. Environmental conditions  
535 within the intestinal tract, such as the presence of bile, bicarbonate, reduced oxygen levels, and  
536 unsaturated fatty acids, play a significant role in promoting the simultaneous expression of genes  
537 responsible for the production of Tcp, CT, and various other genes linked to colonization<sup>12,63</sup>. The  
538 activation of the *ToxT* regulon is also influenced by metabolic cues and quorum sensing<sup>12,63</sup>. Although,  
539 transcription factor binding site prediction algorithms tend to over-predict sites. The correlation of  
540 experimentally determined SNPs with the predicted sites and their different nucleotide frequency  
541 provides a reasonable certainty that the observation reflects the phenomenon. The fact that we found  
542 significant intergenic SNPs in TFBs of 11 TFs and not in promoters, suggests a possible important role  
543 in such scenario. Higher frequency of SNPs close to transcriptional start sites is related to subtle  
544 alteration of gene expression which might result in lineage diversity. In addition to a wider range of  
545 genomic determinants found in this study, we also found 23 genes with mapped SNPs (*tyrA*, *gyrA*, *ctxB*,  
546 *glmM*, *tamA*, *valS*, *czcA*, *licH*, *mutL*, *kbl*, *cobB*, *mak*, *znuC*, *phhA*, *nagA\_1*, *argG*, *cysG\_1*, *murI*, *appC*,  
547 *putA*, *suhB*, *fadJ* and *recD*) in common between our analysis and Monir's comparison of BD-1 vs BD-  
548 2<sup>9</sup> and nine genes with SNPs (*rstA*, *ubiA*, *dsbD*, *clcA*, *thiG*, *rtxA*, *mltD*, *fadJ* and *recD*) in common  
549 between our analysis and Monir's comparison of BD-1.1 vs BD-1.2<sup>10</sup>.

550 Roughly 20% of people who contract toxigenic *V. cholerae* show cholera symptoms<sup>12</sup>. Among  
551 symptomatic cases, approximately 5% are mild, 35% are moderate, and about 60% are severe. The  
552 disease's severity depends on pathogenic factors on the bacteria, and the host, including age, nutrition,

553 and immune system<sup>12</sup>. Here, we revealed the existence of correlations between a core set of genetic  
554 determinants in *V. cholerae* and clinical symptoms (diarrhoeal duration, number of stools, abdominal  
555 pain, vomit, and dehydration). A recent study<sup>65</sup> investigated these correlations, using machine learning,  
556 by analysing gene families in the gut microbiome of household members of Cholera patients to predict  
557 disease severity. In such study, associations were found in gene families like ribosomal proteins, RNA  
558 polymerases, and the sugar phosphotransferase system with symptomatic disease. However, the  
559 computational pipeline adopted in such work<sup>65</sup> did not produce high-performance metrics for predictive  
560 models. Our pipeline, in contrast to Levade *et al*<sup>65</sup>, achieved superior performance metrics, and  
561 encompassed accessory genes, core genome SNPs, and intergenic SNPs. It considered variants in both  
562 functional protein-coding and regulatory forms, revealing their additive effect on diverse clinical  
563 symptoms.

564 Moreover, mechanistic insights were derived through GSMs and protein-protein interaction  
565 networks. Notably, we identified genes crucial for pH homeostasis, host adaptability, colonization,  
566 virulence, motility, acid tolerance, toxin transport, biofilm formation, and bacteriophage resistance.  
567 Important pathways were found underlying these roles, such as the fatty acids biosynthesis which is  
568 important for *V. cholerae* since unsaturated fatty acids present in bile inhibit the expression of virulence  
569 factors and both cholesterol and unsaturated fatty acids can enhance the motility of *V. cholerae*<sup>66</sup>; and  
570 biofilm production which plays a crucial role in the cholera pathogenesis and dissemination of disease<sup>62</sup>.  
571 Furthermore, our ML analysis identified genes associated to abdominal pain that were also found  
572 important for colonization in *V. cholerae*. It is known that colonization of pathogenic bacteria can  
573 present clinical symptoms such as abdominal pain<sup>67</sup>.

574 Three non-synonymous SNPs associated to the clinical symptoms were also found as statistically  
575 significant in differentiating the BD-1.2 and BD-2 lineages. These SNPs mapped to *clpS*, *gshB* and  
576 *fabV*. In *V. cholerae* ClpS regulation involves cAMP receptor protein (CRP)<sup>35</sup>. CRP is important in *V.*  
577 *cholerae* gene regulatory network lifestyle switching, adapting gene expression for quorum sensing,  
578 intestinal colonization, and toxin production to its environment<sup>35</sup>. GshB, encodes a glutathione  
579 synthetase (GSH), a gene associated to resistance to oxidative stress. It is part of the  $\sigma$ 32 regulon,

580 contributing to *V. cholerae* intestinal colonization<sup>37</sup>. Glutathione controls the potassium efflux system,  
581 Kef, and pH homeostasis involved in Na<sup>+</sup> and K<sup>+</sup> transport<sup>68</sup>. Impaired glutathione production may  
582 affect the stress response<sup>68</sup>. *GshB* was additionally shown to have a role in the ability of *V. cholerae* to  
583 mount an acid tolerance response<sup>38</sup>. *V. cholerae fabV* is one of the several triclosan-resistant ENR  
584 encoding genes<sup>36</sup>. Resistance to triclosan also affects resistance to other antibiotics, showing cross-  
585 resistance to a wide range of antibiotics (including chloramphenicol and tetracycline)<sup>69</sup>. Moreover, *fabV*  
586 exhibits pleiotropic effects controlling pathogenicity in *P. aeruginosa* via modulation of fatty acids  
587 synthesis, production of virulence factors and motility<sup>70</sup>.

588 Analysing the 3D structure based on non-synonymous mutations can provide insights into the  
589 mechanisms by which these mutations can cause disease<sup>71-74</sup>. Changes in the stability of proteins can  
590 lead to manifestation of diseases<sup>73</sup> or symptom variations<sup>71,74</sup>. Among all types of mutations, non-  
591 synonymous SNPs have the greatest impact on protein structure and function<sup>75</sup>. In this work we found  
592 that different SNPs accumulated in BD-1.2 isolates compared to BD-2 isolates, suggesting different  
593 evolutionary dynamics possibly explaining the temporal shift of the two lineages. Our analysis of top-  
594 ranked non-synonymous SNPs in protein-coding regions, identified by machine learning as linked to  
595 both BD-1.2 lineage evolution and clinical symptoms, specifically FabV and GshB, unveiled that SNPs  
596 present in BD-1.2, associated with more severe cholera, led to increased protein stability. That protein  
597 stability might be relevant for disease severity is also supported by the fact that no SNPs associated to  
598 clinical symptoms were found in any TFBS or promoter signature but only in protein-coding sequences.  
599 In this study, we have identified promising targets related to metabolism (*clcA*, *cysG*, *adh*),  
600 antimicrobial resistance (i.e. *bcr*, *blc*), and virulence (i.e. *ompU*, *skp*, *tamA*, *vals*). These targets show  
601 significant potential for further investigation through experimental studies.

602 We are aware of the limitations of our current study. Several host factors (retinol deficiency, blood  
603 group, genetic factors, innate immune system) confer susceptibility to cholera with higher risk of  
604 symptomatic disease<sup>76</sup>. These factors have not been considered in this study due to lack of data. A  
605 further limitation of this study was the inability to consider the potential impact of co-infections with  
606 either multiple *V. cholerae* lineages/strains or other pathogens. Whilst the presence of more than one *V.*



607 *cholerae* strain or lineage in a host has recently been shown to be unlikely<sup>77-79</sup>, co-infections with other  
608 bacteria can occur in diarrheal patients. A study of 10,351 confirmed clinical *V. cholerae* cases from  
609 2000-2021 in Bangladesh found that *Campylobacter* spp., enterotoxigenic *E. coli* (ETEC) and rotavirus  
610 were the most frequently found co-pathogens, with co-infection rates of 6.7%, 5.7% and 2.4%  
611 respectively<sup>80</sup>. Although the effects on the host of co-infection of *V. cholerae* with *Campylobacter* spp.  
612 or rotavirus have not been studied, co-infection with enterotoxigenic *E. coli* (ETEC) has been studied.  
613 Chowdhury *et al* 2010<sup>81</sup> showed that coinfection with ETEC results in an increased host immune  
614 response, and so could potentially affect observed symptoms. The authors have also observed a higher  
615 co-infection rate (13%) between *V. cholerae* O1 and ETEC in their cohort. However, for future research  
616 will aim to incorporate these variables to provide a more comprehensive understanding of the  
617 interactions between host and pathogen, as well as between different pathogens, in the context of  
618 cholera. This study should be considered a proof-of-principle to be further investigated and validated  
619 with larger sample sizes and different geographical areas. With the advent of modern technologies, by  
620 strengthening bespoke analytical methods and by performing wider comparisons (asymptomatic vs.  
621 symptomatic, patients vs. households, environmental vs stool vibrio) we can potentially disentangle the  
622 intricate network of correlations between the genetic underpinnings of cholera symptoms and  
623 epidemiological transmission risk, uncovering regulatory, metabolic and signalling networks  
624 interconnectivity that might help to inform future interventions.

625

## 626 **Methods**

### 627 **Ethics Statement**

628 Informed written consent was obtained from all adult patients, or guardians on behalf of children. Upon  
629 receiving consent, the physician collected the patient's sociodemographic characteristics and medical  
630 histories. For the icddr,b isolates, the study protocol was approved by the Institutional Review Board  
631 of icddr,b (PR-15127). For the IEDCR isolates, the study was performed in accordance with protocols

632 approved by the Institutional review board of IEDCR (IEDCR/IRB/09 and IEDCR/IRB/26). Ethics  
633 approval was also obtained from the University of Nottingham (2811 110724).

### 634 **Experimental Design**

635 For the study we used 129 *V. cholerae* bacterial isolates obtained from distinct stool samples of patients  
636 between 2014 and 2021 from the ongoing Nationwide Cholera Surveillance<sup>82</sup>, jointly conducted by  
637 IEDCR and icddr,b. The isolates were collected from admitted patients from six divisions of Bangladesh  
638 (Barisal n=11, Chittagong n=6, Dhaka n=99, Khulna n=2, Rajshahi n=4, and Sylhet n=7). The isolates  
639 included in the study were gathered from patients meeting the case definition of diarrhoea and  
640 consenting to be included in the surveillance study. The case definition was used and defined as: i)  
641 Diarrhoea (patient age > 2 months): any patient attending hospital with 3 or more loose or liquid stools  
642 within 24 hours or less than 3 loose / liquid stools causing dehydration; ii) Diarrhoea (patient age < 2  
643 months): changed stool habit from usual pattern in terms of frequency (more than the usual number of  
644 purging) or nature of stool (more water than faecal matter). The case definition of diarrhoea was  
645 standardized to ensure consistency across different regions and over the collection timeline. Stool  
646 samples were processed by either IEDCR or icddr,b research institutes. For the identification of *V.*  
647 *cholerae*, specimens were streaked onto taurocholate-tellurite gelatin agar (TTGA) and incubated  
648 overnight at 37°C. Specimens were also inoculated in alkaline peptone water for enrichment and  
649 incubated for an additional 18–24 hours<sup>83</sup> and plated on TTGA. Suspected colonies were serotyped with  
650 monoclonal antibody specific to *V. cholerae* O1 (Ogawa and Inaba) and O139 serogroups<sup>84</sup> for the  
651 icddr,b isolates, while for the IEDCR isolates serotyping and biotyping was carried out by slide  
652 agglutination and PCR using primers in Supplementary Data 21. Further confirmation of the isolates  
653 being *V. cholerae* was obtained by whole genome sequencing. Confirmed isolates were tested for  
654 antimicrobial susceptibility using disk diffusion methods in accordance with CLSI protocols<sup>85</sup> to  
655 antibiotics: ampicillin, azithromycin, ciprofloxacin, ceftriaxone, cefixime, doxycycline, erythromycin  
656 and meropenem, using commercially available antibiotic discs (Oxoid, Basing- stoke, United  
657 Kingdom). *Escherichia coli* American Type Culture Collection 25922 susceptible to all antimicrobials  
658 was used as a control strain for susceptibility studies.

659 Clinical metadata was collected from patients corresponding to 104 isolates for the 129 isolates in our  
660 cohort. Clinical data covered 5 categories (duration of diarrhoea, number of stools, abdominal pain,  
661 vomiting, and dehydration), in addition the age and sex of the patient and location of the patient was  
662 recorded. Clinical symptoms data (Supplementary Data 14) were binned into categories and ranked in  
663 order of increasing severity for data analysis.

664 • Duration of diarrhoea: number of days the diarrhoea persisted was recorded. Data were binned as a  
665 duration score ranging from 1-3, with 1 = <1 day; 2 = 1-3 days; 3 = 4-6 days.

666 • Number of stools in 24 hours: The number of stools recorded in a 24-hour period during the hospital  
667 admission was recorded. Data were binned as a number of stools score ranging from 1-5 with 1= 3-5  
668 times; 2= 6-10 times; 3=11-15 times; 4=16-20 times; 5=21+ times.

669 • Abdominal pain: the presence or absence of abdominal pain was recorded as a 0 for absence and 1 for  
670 present.

671 • Vomit: The presence or absence of any vomiting in the 24 hours prior to admission was recorded with  
672 0 denoting no vomiting and 1 denoting the occurrence of vomiting

673 • Dehydration: clinical assessment of dehydration was recorded as none, moderate or severe by the  
674 clinician.

### 675 **DNA purification and extraction**

676 DNA extraction was performed at North South University. All the *V. cholerae* isolates were subjected  
677 to genomic DNA extraction in accordance with the manufacturers protocol of QIAamp DNA Mini Kit  
678 (Qiagen).

### 679 **Library construction and whole-genome sequencing**

680 The library preparation and sequencing of the 129 selected strains were carried out at NGRI (NSU  
681 Genomics Research Institute, North South University). To prepare the Illumina libraries, approximately  
682 1 µg of high molecular weight *V. cholerae* genomic DNA was utilized. Barcoded libraries were prepared  
683 using the Illumina DNA Prep Kit (product code 20060059, NEB, USA) following the manufacturers  
684 protocol. Nextera DNA CD index codes were added to attribute sequences to each sample. Following

685 that, paired-end sequencing with  $2 \times 151$  cycles was performed on the Illumina MiSeq platform at  
686 NGRI.

### 687 **Genome assembly and annotation**

688 All sequences were pre-processed to using the Illumina BaseSpace sequencing hub. To clean the data  
689 adapters were trimmed and unidentified bases were removed. Genomes were assembled using SPAdes  
690 (v3.12)<sup>86</sup> with default parameters and a coverage cut off value of 20. Genomic contamination was  
691 assessed using ContEst16S<sup>87</sup> with only genomes identified as *V. cholerae* retained for further analysis.  
692 Contigs with length shorter than 500 nucleotides were filtered out of the final assemblies. Genomes  
693 were annotated with Prokka (v1.14.6)<sup>88</sup>, using default settings with `--addgenesz--usegenus`.

694 Screening of annotated genes against ABR databases, virulence and plasmid databases and in silico  
695 subtyping.

696 The whole-genome sequences were screened against the CARD<sup>89</sup> database (accessed 05-06-2022) using  
697 Abricate<sup>90</sup> with a minimum coverage of 70% and minimum identity of 90% to identify known AMR-  
698 associated genes in the isolate cohort. Genomes were also screened against the VFDB<sup>91</sup> database using  
699 Abricate<sup>90</sup> to find virulence associated genes, with 70% coverage and 90% identity) (accessed 05-06-  
700 2022). Plasmids screening was conducted using the PlasmidFinder<sup>92</sup> database in Abricate<sup>90</sup>, with 70%  
701 coverage and 90% identity) (accessed 05-06-2022); no plasmids were identified in the genome  
702 sequences. Sequence types were identified through MLST<sup>93</sup> which mapped the sequences to the  
703 PubMLST<sup>94</sup> database.

### 704 **Pangenome analysis and generation of genetic features input files**

705 All annotated genomes we used as input for pangenome analysis using Roary v3.13<sup>95</sup>. The core genome  
706 alignment was taken as input to produce a file of core gene SNPs present in the cohort using SNP sites  
707 2.5.1<sup>96</sup>. SNPs within intergenic regions (IGRs) were extracted using piggy v1.5<sup>97</sup> to generate an  
708 alignment of core intergenic clusters. Variants in this alignment were then called using SNP sites 2.5.1.  
709 The presence-absence of accessory gene was found from the output of Roary.

710 In addition, a further pangenome alignment was created consisting of the 129 isolates in our cohort  
711 together with 218 isolates collected in Bangladesh from 2004 to 2022 (The European Nucleotide  
712 Archive-ENA (<http://www.ebi.ac.uk/ena>), accession codes: PRJDB8664, PRJDB12727, PRJDB13928,  
713 PRJNA723557).

#### 714 **Phylogenetic analysis of *V. cholerae* isolates in our cohort in Bangladesh**

715 For both our cohort alone and our cohort together with publicly available Bangladeshi isolates (as  
716 detailed above) maximum likelihood phylogenies were reconstructed. Using the core genome  
717 alignments generated in Roary v3.13<sup>95</sup>, the phylogenies were reconstructed in IQ Tree (v2.2.0.3)<sup>98</sup> with  
718 10000 ultrafast bootstrap replicates and best fitted evolutionary model (HKY+F+I for our cohort only  
719 and K3Pu+F+I for the combined Bangladesh alignment) was selected using ModelFinder<sup>99</sup>. The  
720 alignment length of the core genome of our cohort was 3459819 nucleotide sites of which 1486 were  
721 informative. For the core genome of the combined Bangladeshi isolates, the alignment length was  
722 2086397 nucleotide sites with 844 informative sites. The resulting consensus trees were visualised using  
723 iTol v6<sup>100</sup>, and branches with less than 95% ultrafast bootstrap support were deleted.

#### 724 **Phylogenetic relations between *V. cholerae* isolates worldwide**

725 We used WGS data from 1140 *V. cholerae* isolates collected from India, Africa, Haiti, Yemen together  
726 with our Bangladesh samples (see Supplementary Data 2 and 3). To generate the input for a  
727 phylogenetic tree, SNP variants were called from each isolate against the reference genome VC N16961  
728 (NC\_002505.1; NC\_002506.1) using Snippy v4.6.0<sup>101</sup> (<https://github.com/tseemann/snippy>). The  
729 cleaned alignment files from Snippy were concatenated via the SeqIO function of biopython v1.83<sup>102</sup>  
730 then recombination was masked using Gubbins (v.2.3.4)<sup>103</sup>. The filtered polymorphic sites output from  
731 Gubbins was further filtered using SNP-sites<sup>96</sup>. The final SNP input contained 4033464 nucleotide sites  
732 with 26995 informative sites. This recombination-free SNP output was then used as input to reconstruct  
733 the phylogeny using IQtree (v2.2.0.3)<sup>98</sup> with 1000 ultrafast bootstrap replicates and best fitted model  
734 (K3Pu+F+I+G4) was selected by ModelFinder<sup>99</sup>. The sequence ERR025382 (Indonesia-1957) was used  
735 as an outgroup, and the tree was rooted here. The resulting consensus tree was visualised using iTol  
736 v6<sup>93</sup>, and branches with less than 95% ultrafast bootstrap support were deleted.

737 **Transcriptional binding motifs**

738 Motif searches were conducted using FIMO (Find Individual Motif Occurrences<sup>104</sup> within the MEME  
739 (Multiple Em for Motif Elicitation)<sup>105</sup> suite (<https://meme-suite.org/meme/tools/fimo>). Reference  
740 sequences of intergenic regions of DNA from our isolates were generated in Piggy as described above;  
741 these were used as input for FIMO. To predict the TFBS the following databases were used: CollecTF  
742 (Bacterial TF Motifs); Prokaryotes (Prodoric Release 8.9); Prokaryotes (RegTransBase v4); Combined  
743 Prokaryotes. Intergenic regions where motifs were found were variant called using SNP-sites<sup>96</sup> and then  
744 aligned to the motif sequences using Clustal Omega v1.2.4<sup>106</sup>. For visualisation of intergenic regions,  
745 alignment maps of the intergenic regions were created using Jalview 2.11.3.2 with easyfig python  
746 genome figure package<sup>107</sup>.

747 **Promoter analysis for Intergenic SNPs**

748 BPRM/softberry<sup>108</sup> was used to predict promoter region and oligonucleotides from known TF binding  
749 sites close to the promoter region.

750 **Genome-scale metabolic model**

751 All simulations were performed using the Python cobra toolkit v0.26.2. The analysis was conducted on  
752 both a manually curated and validated model of *V. cholerae* O1 N16961, iAM-Vc960, taken from  
753 Abdel-Haleem *et al*<sup>19</sup> and on automatically generated draft strain-specific GSM models. The strain-  
754 specific draft models were generated using CarveMe<sup>27</sup>. CarveMe was run using the CPLEX solver and  
755 gram negative template, with gap filling for LB and M9 media using the command: ‘carve input.faa --  
756 gapfill M9,LB -u gramneg --solver cplex --output model.xml’. Gene essentiality, FVA and FBA  
757 analyses as described below were conducted on genes of interest in the generalised iAM-Vc960 and in  
758 each of the 129 draft strain-specific models, based on the analysis pipeline in Pearcy *et al*<sup>39</sup>.

759 For all gene essentiality, FBA and FVA analyses, a knockout model for each gene of interest was  
760 constructed by blocking all corresponding reactions to zero, given that the reaction is not catalysed by  
761 an isozyme. We considered the essentiality of a gene under both rich medium conditions and M9  
762 minimal medium conditions. To mimic rich medium conditions, the model was constrained to allow all  
763 carbon sources into the system, with a fixed uptake rate of 1 mmol/gDCW/h. If a feasible solution exists,

764 while maximizing the biomass equation as the objective function, then the knockout of the gene was  
765 not essential. To mimic M9 minimal medium conditions, the model was constrained so one individual  
766 carbon source had a maximum uptake of 10 mmol/gDCW/h. This simulation (minimal medium  
767 condition) was repeated for each carbon source in the model. The genes whose corresponding knockout  
768 model achieved a growth rate of 0.0001 h<sup>-1</sup> or less were considered essential. Flux variability analysis  
769 (FVA) was applied to the wild-type model and each knockout model using the cobra toolbox in  
770 python<sup>109</sup>. FVA calculates the minimum and maximum flux through each reaction in the model, given  
771 a set of constraints, resulting in the range of possible fluxes for each reaction (flux span). FVA was  
772 simulated using glucose as the only carbon source in aerobic minimal M9 medium conditions. Note that  
773 reaction loops in the solution were not allowed. A gene knockout was considered to significantly affect  
774 the flux if the flux span of at least one reaction was changed by greater than 10% compared to the  
775 wildtype solution. For the FBA analysis, a drain reaction (i.e., a reaction that consumes the metabolite  
776 of interest) was added to the GSM model for each metabolite. The maximum theoretical yield of each  
777 metabolite was calculated by setting its corresponding drain reaction as the objective function, with  
778 glucose as the only carbon source in aerobic minimal M9 medium conditions. All metabolites contained  
779 within the model were considered in the FBA analysis. In iAM-Vc960 this was 1,741 different  
780 metabolites, whilst in the draft strain-specific GSM models the number of metabolites spanned the range  
781 1321-1433. The simulations were carried out for the wild-type model and each gene knockout model.  
782 A gene knockout was considered to significantly affect metabolite yield if the yield of at least one  
783 metabolite was reduced to zero, given that it was non-zero in the wildtype. For each of the selected  
784 genes of interest, molecular function, pathways and biological processes were taken from the BioCyc  
785 database<sup>110</sup> using the SMART tables for *Vibrio cholerae* O1 biovar El Tor strain N16961. These were  
786 added to Supplementary Data 11, 12, 17 and 18 to give context to the analysis results.

### 787 **Network analysis based on core genome SNPs**

788 Network of our cohort of 129 *V. cholerae* isolates was created using a pairwise hamming distance  
789 comparison based on core genome SNPs in python (NetworkX v2.8.4<sup>111</sup> and Matplotlib v3.6.2<sup>112</sup>). Each  
790 node represents an isolate while the edge represents the hamming distance between two isolates

791 multiplied by the total number of SNPs found in our cohorts (2,382 SNPs). A threshold of 15 or less  
792 SNPs difference was used to filter the edges in the network as suggested by Ludden et al (2019)<sup>113</sup> and  
793 used by us previously<sup>17,18</sup>.

#### 794 **Statistical analysis and machine learning of genomic features correlated to a specific lineage or** 795 **clinical symptoms**

796 To assess if the genomic features were associated with a lineage or to a clinical symptom, we employed  
797 a fisher exact test<sup>9,10</sup>. Furthermore, to analyse the relationship between genomic features of the BD-1.2  
798 lineage and clinical symptoms a machine learning pipeline was employed. Clinical data were collected  
799 from 104 out of 129 *V. cholerae* isolates of which 63 belonged to the BD-1.2 lineage. These clinical  
800 symptoms were be divided into two groups: binary (vomit and abdominal pain) and multi-class  
801 (dehydration, number of stools and duration of diarrhoea), with the binning within each group described  
802 above. In the multiclass group, we applied a one-vs-one approach, i.e., each class is compared  
803 individually to another class. For example, dehydration class Moderate is compared against class  
804 Severe. For both binary and multiclass groupings, as the classes were unbalanced, we oversampled the  
805 minority class as a pre-processing step using a Synthetic Minority Over-sampling Technique approach  
806 (SMOTE)<sup>114</sup>. The Python package Scikit-learn version 1.2.1<sup>115</sup> was used to make the classification and  
807 the package Scipy version 1.9.3<sup>116,117</sup> was used to select the most important features based on a Fisher  
808 exact test.

809 The pipeline first removes features that are either present or absent in all the samples. Second, to  
810 measure the influence of confounding effects in the data, it uses a two-sided chi-square test of  
811 independence to measure the dependency between the confounding effects (sex of patient, age of  
812 patients, year of collection, location of patient, serology of *V. cholerae*) and the phenotype classes (p-  
813 value < 0.01 with Bonferroni correction); if the null hypothesis is rejected (i.e. there is a dependency  
814 between the confound effect and the phenotype) the pipeline checks if there are features that are  
815 dependent on the confounding effect again based on a two-sided chi-square test of independence (p-  
816 value < 0.01 with Bonferroni correction); if there are features where the null hypothesis is rejected,  
817 these features are removed from the analysis. Next, the pipeline oversamples the minority class using a



818 SMOTE approach. Then, based on the oversampled data, it selects the most important features using a  
819 two-sided Fisher exact test ( $p\text{-value} < 0.1$ ). This process is done in two parts: i) to improve  
820 randomization in the pipeline and avoid confounding effects, a loop over 1000 different random seeds  
821 is used for the SMOTE approach in order for it to have different initializations; for each loop the most  
822 important features are selected based on the Fisher exact test; ii) then, the features that are selected in  
823 over 75% of the different initializations are deemed important and a random initialization is selected  
824 that contains all these important features to be used for the prediction models. Next, a panel of machine  
825 learning methods (logistic regression (LR), linear support vector machine (L-SVM), radial basis  
826 function support vector machine (RBF-SVM), extra-tree classifier, random forest, Adaboost and  
827 XGboost) was used to predict the clinical symptoms classes based on the pre-selected features described  
828 above. The hyperparameters used were:

- 829 • Logistic Regression: inverse of regularization strength  $C = [0.0001, 0.001, 0.01, 0.1, 1, 10, 100,$   
830  $1000, 10000]$ ;
- 831 • Linear SVM: penalty parameter of the hinge loss error  $C = [0.0001, 0.001, 0.01, 0.1, 1, 10, 100,$   
832  $1000, 10000]$ ;
- 833 • Random Forests, Extra Trees and Adaboost: Number of estimators =  $[2, 4, 8, 16, 32, 64, 128,$   
834  $256]$ ;
- 835 • Non-linear SVM with RBF kernel:  $\gamma$  (RBF kernel coefficient) =  $[0.0001, 0.0001, 0.001, 0.01,$   
836  $0.1, 1]$  and  $C$  (L2 penalty parameter) =  $[0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]$ ;
- 837 • XGBoost: Number of estimators =  $[2, 4, 8, 16, 32, 64, 128, 256]$  and learning rate =  $[0.0001,$   
838  $0.001, 0.01, 0.1, 1]$ .

839 As per previous works<sup>17,18,29,30</sup>: (i) nested cross-validation<sup>118,119</sup> was employed to assess the performance  
840 and select the hyper-parameters of the proposed classifiers and to compare the results obtained by the  
841 seven different classifiers used; (ii) a Friedman Statistical F-test ( $F_F$ ) with Iman-Davenport correction  
842 was used for statistical comparison of multiple classifiers across multiple analyses<sup>120</sup>; (iii) a post-hoc  
843 Nemenyi test was employed to find if there is a single classifier or a group of classifiers that performs  
844 statistically better in terms of their average AUC rank after the  $F_F$  test has rejected the null hypothesis

845 (stating that the performance of the comparisons on the individual classifiers over the different datasets  
846 is similar)<sup>120</sup>; (iv) an undirected graph was created using NetworkX<sup>111</sup> to visualize how the features  
847 (accessory genes, core genome SNPs and intergenic SNPs) correlate with different clinical symptoms.

### 848 **Protein-protein interaction network and building protein 3D structures**

849 Protein-protein interaction networks of the protein encoded of the genes associated with clinical  
850 symptoms were obtained using STRING database v12.0 (using reference genome *V. cholerae* O1 biovar  
851 El Tor str. N16961) and analysed in Cytoscape 3.10.1<sup>121</sup>. Eighty-one accessory and core genes selected  
852 by machine learning were used as input for the PPI, of these only 60 could be mapped to the STRING  
853 database. The interactome was constructed using first and second neighbour proteins. Disconnected  
854 nodes and nodes with interaction scores lower than medium confidence level (interaction scores  
855 <0.400), according to StringDB<sup>122</sup>, were filtered out. Functions of the protein in the network were  
856 annotated with Gene Ontology terms (biological process, molecular function, cellular component and  
857 KEGG pathways) in StringDB<sup>122</sup>. Three-dimensional AlphaFold<sup>123</sup> predicted models were obtained by  
858 aligning the protein FASTA sequence to reference sequences from the Uniprot database<sup>124</sup> to find a 3D  
859 protein structure. 3D protein structures were then visualised using UCSF Chimera<sup>125</sup> and UCSF  
860 ChimeraX<sup>126</sup>. Protein stability analysis and the effect of each mutation were performed with DUET<sup>127</sup>,  
861 DynaMut<sup>128</sup> and SIFT<sup>129</sup>. The electrostatic potential was analysed and visualised using PDB2PQR and  
862 APBSaccessed online<sup>130</sup>, UCSF ChimeraX<sup>126</sup> and APBS Coloring<sup>130</sup>.

### 863 **Statistical Analysis**

864 Statistical comparisons were made using the SciPy package implementing: 1. A two-sided chi-squared  
865 test with Bonferroni correction to evaluate the similarities between the serotypes and the collection year  
866 and location of the isolates (p-value < 0.005); 2. A two-sided Mann Whitney U test to evaluate the  
867 distribution of the counts of accessory genes, coding and non-coding SNPs in BD-1.2 and BD-2 lineages  
868 and along the different collection years (p value < 0.005); 3. A two-sided Fisher exact test, with  
869 Bonferroni correction, to assess the relationship between the BD-2 and BD-1.2 lineages and different  
870 genomic features - core and intergenic SNPs and accessory genes (p value < 0.005); 4. A two-sided  
871 hypergeometric enrichment tests (two-sided) with false discovery rate (FDR) for the GSM analysis (p-

872 value < 0.01); and 5. A two-sided chi-square test of independence to test if there are symptoms/features  
873 that are dependent on the confounding effect (p-value 0.01 with Bonferroni correction); 6. A two-sided  
874 Fisher exact test to select the most important features in the machine learning pipeline (p-value < 0.1);  
875 7. A two-sided Friedman Statistical F-test (FF) with Iman-Davenport correction for statistical  
876 comparison of multiple datasets over the seven different classifiers used (p-value < 0.05). With 7  
877 classifiers and 6 clinical symptom models, the Friedman test is distributed according to the F  
878 distribution with  $7-1 = 6$  and  $(7-1) \times (6-1) = 30$  degrees of freedom. Therefore, the critical values for  
879  $F(6,30)$  using a p-value = 0.05 is 2.42052319. The post-hoc Nemenyi test was used to find if there is a  
880 single classifier or a group of classifiers that performs statistically better in terms of their average rank  
881 after the FF test has rejected the null hypothesis (stating that the performance of the comparisons on the  
882 individual classifiers over the different datasets is similar); 8. A two-sided Mann Whitney U test was  
883 used to assess for lineage differences in the numbers of genes, reactions and metabolites in the generated  
884 draft strain-specific GSM models. 9. A two-sided Mann Whitney U test was used to assess the number  
885 of affected reactions and metabolites in knockouts of genes discriminating lineages, compared to  
886 randomly selected genes.

## 887 **Data Availability**

888 Short-read sequence data for all 129 isolates used in this study are deposited in the NCBI SRA and can  
889 be found associated with BioProject number PRJNA1021874  
890 [<https://www.ncbi.nlm.nih.gov/bioproject/?term=PRJNA1021874>]. All previously published public *V.*  
891 *cholerae* sequences used in this study are held in European Nucleotide Archive-ENA or NCBI  
892 repositories under accession numbers supplied in Supplementary Data 2. Reference sequences are  
893 available from NCBI under accessions: NC\_002505.1  
894 [[https://www.ncbi.nlm.nih.gov/nuccore/NC\\_002505.1](https://www.ncbi.nlm.nih.gov/nuccore/NC_002505.1)], NC\_002506.1  
895 [[https://www.ncbi.nlm.nih.gov/nuccore/NC\\_002506.1](https://www.ncbi.nlm.nih.gov/nuccore/NC_002506.1)] and European Nucleotide Archive-ENA under  
896 accession: ERR025382 [<https://www.ebi.ac.uk/ena/browser/view/ERR025382>]. Clinical data used in  
897 this study is given in Supplementary Data 14.

## 898 Code Availability

899 The code used in this study and draft strain-specific GSMMs are available in the following GitHub  
900 repository: <https://github.com/tan0101/VibrioCARE> under  
901 <https://doi.org/10.5281/zenodo.13325928><sup>131</sup>.

902

## 903 References

904

- 905 1 Baddam, R. *et al.* Genome dynamics of *Vibrio cholerae* isolates linked to seasonal outbreaks  
906 of cholera in Dhaka, Bangladesh. *MBio* **11**, e03339-03319 (2020).
- 907 2 Banerjee, R., Das, B., Nair, G. B. & Basak, S. Dynamics in genome evolution of *Vibrio*  
908 *cholerae*. *Infect. Genet. Evol.* **23**, 32-41 (2014).
- 909 3 Ali, M., Nelson, A. R., Lopez, A. L. & Sack, D. A. Updated Global Burden of Cholera in  
910 Endemic Countries. *PLoS Negl. Trop. Dis.* **9**, e0003832, doi:10.1371/journal.pntd.0003832  
911 (2015).
- 912 4 Kaper, J. B., Morris, J. G., Jr. & Levine, M. M. Cholera. *Clin. Microbiol. Rev.* **8**, 48-86,  
913 doi:10.1128/cmr.8.1.48 (1995).
- 914 5 Karaolis, D. K. *et al.* A *Vibrio cholerae* pathogenicity island associated with epidemic and  
915 pandemic strains. *Proceedings of the National Academy of Sciences* **95**, 3134-3139 (1998).
- 916 6 Son, M. S., Megli, C. J., Kovacicova, G., Qadri, F. & Taylor, R. K. Characterization of *Vibrio*  
917 *cholerae* O1 El Tor biotype variant clinical isolates from Bangladesh and Haiti, including a  
918 molecular genetic analysis of virulence genes. *J. Clin. Microbiol.* **49**, 3739-3749 (2011).
- 919 7 Wozniak, R. A. *et al.* Comparative ICE genomics: insights into the evolution of the  
920 SXT/R391 family of ICEs. *PLoS Genet.* **5**, e1000786, doi:10.1371/journal.pgen.1000786  
921 (2009).
- 922 8 Faruque, S. M. & Mekalanos, J. J. Pathogenicity islands and phages in *Vibrio*  
923 *cholerae* evolution. *Trends Microbiol.* **11**, 505-510, doi:10.1016/j.tim.2003.09.003  
924 (2003).
- 925 9 Monir, M. M. *et al.* Genomic Characteristics of Recently Recognized *Vibrio cholerae* El Tor  
926 Lineages Associated with Cholera in Bangladesh, 1991 to 2017. *Microbiology Spectrum* **10**,  
927 e00391-00322 (2022).
- 928 10 Monir, M. M. *et al.* Genomic attributes of *Vibrio cholerae* O1 responsible for 2022 massive  
929 cholera outbreak in Bangladesh. *Nat. Commun.* **14**, 1154, doi:10.1038/s41467-023-36687-7  
930 (2023).
- 931 11 Morita, D. *et al.* Whole-genome analysis of clinical *Vibrio cholerae* O1 in Kolkata, India, and  
932 Dhaka, Bangladesh, reveals two lineages of circulating strains, indicating variation in  
933 genomic attributes. *Mbio* **11**, e01227-01220 (2020).
- 934 12 Baker-Austin, C. *et al.* *Vibrio* spp. infections. *Nat Rev Dis Primers* **4**, 8, doi:10.1038/s41572-  
935 018-0005-8 (2018).
- 936 13 Domman, D. *et al.* Integrated view of *Vibrio cholerae* in the Americas. *Science* **358**, 789-793,  
937 doi:10.1126/science.aao2136 (2017).
- 938 14 Rashid, M. U. *et al.* CtxB1 outcompetes CtxB7 in *Vibrio cholerae* O1, Bangladesh. *J. Med.*  
939 *Microbiol.* **65**, 101-103, doi:10.1099/jmm.0.000190 (2016).
- 940 15 Jubya, F. T. *et al.* *Vibrio cholerae* O1 associated with recent endemic cholera shows  
941 temporal changes in serotype, genotype, and drug-resistance patterns in Bangladesh. *Gut*  
942 *Pathog.* **15**, 17, doi:10.1186/s13099-023-00537-0 (2023).

- 943 16 Weill, F. X. *et al.* Genomic history of the seventh pandemic of cholera in Africa. *Science* **358**,  
944 785-789, doi:10.1126/science.aad5901 (2017).
- 945 17 Baker, M. *et al.* Convergence of resistance and evolutionary responses in *Escherichia coli* and  
946 *Salmonella enterica* co-inhabiting chicken farms in China. *Nat. Commun.* **15**, 206,  
947 doi:<https://doi.org/10.1038/s41467-023-44272-1> (2024).
- 948 18 Peng, Z. *et al.* Whole-genome sequencing and gene sharing network analysis powered by  
949 machine learning identifies antibiotic resistance sharing between animals, humans and  
950 environment in livestock farming. *PLoS Comput. Biol.* **18**, e1010018, doi:  
951 <https://doi.org/10.1371/journal.pcbi.1010018> (2022).
- 952 19 Abdel-Haleem, A. M. *et al.* Integrated metabolic modeling, culturing, and transcriptomics  
953 explain enhanced virulence of *Vibrio cholerae* during coinfection with enterotoxigenic  
954 *Escherichia coli*. *MSystems* **5**, e00491-00420 (2020).
- 955 20 Karp, P. D. *et al.* The ecocyc database. *EcoSal Plus* **8** (2018).
- 956 21 Zhang, H., Luo, Q., Gao, H. & Feng, Y. A new regulatory mechanism for bacterial lipolic acid  
957 synthesis. *MicrobiologyOpen* **4**, 282-300 (2015).
- 958 22 Ramamurthy, T. *et al.* Virulence regulation and innate host response in the pathogenicity of  
959 *Vibrio cholerae*. *Frontiers in Cellular and Infection Microbiology* **10**, 572096 (2020).
- 960 23 Jugder, B.-E. *et al.* *Vibrio cholerae* high cell density quorum sensing activates the host  
961 intestinal innate immune response. *Cell Rep.* **40**, 111368 (2022).
- 962 24 Somerville, G. A. *et al.* Correlation of Acetate Catabolism and Growth Yield in  
963 *Staphylococcus aureus*: Implications for Host-Pathogen Interactions. *Infect. Immun.*  
964 **71**, 4724-4732, doi:doi:10.1128/iai.71.8.4724-4732.2003 (2003).
- 965 25 Keating, T. A., Marshall, C. G. & Walsh, C. T. Vibriobactin biosynthesis in *Vibrio cholerae*:  
966 VibH is an amide synthase homologous to nonribosomal peptide synthetase condensation  
967 domains. *Biochemistry* **39**, 15513-15521, doi:10.1021/bi001651a (2000).
- 968 26 Kostiuik, B. *et al.* *Vibrio cholerae* Alkalizes Its Environment via Citrate Metabolism to Inhibit  
969 Enteric Growth *In Vitro*. *Microbiology Spectrum* **11**, e04917-04922,  
970 doi:doi:10.1128/spectrum.04917-22 (2023).
- 971 27 Machado, D., Andrejev, S., Tramontano, M. & Patil, K. R. Fast automated reconstruction of  
972 genome-scale metabolic models for microbial species and communities. *Nucleic Acids Res.*  
973 **46**, 7542-7553, doi:10.1093/nar/gky537 (2018).
- 974 28 Cakar, F., Zingl, F. G., Moisi, M., Reidl, J. & Schild, S. In vivo repressed genes of *Vibrio*  
975 *cholerae* reveal inverse requirements of an H<sup>+</sup>/Cl<sup>-</sup> transporter along the gastrointestinal passage. *Proceedings of the National Academy of*  
976 *Sciences* **115**, E2376-E2385, doi:doi:10.1073/pnas.1716973115 (2018).
- 977 29 Baker, M. *et al.* Machine learning and metagenomics reveal shared antimicrobial resistance  
978 profiles across multiple chicken farms and abattoirs in China. *Nature Food* **4**, 707-720,  
979 doi:10.1038/s43016-023-00814-w (2023).
- 980 30 Maciel-Guerra, A. *et al.* Dissecting microbial communities and resistomes for interconnected  
981 humans, soil, and livestock. *The ISME Journal* **17**, 21-35, doi:[https://doi.org/10.1038/s41396-](https://doi.org/10.1038/s41396-022-01315-7)  
982 [022-01315-7](https://doi.org/10.1038/s41396-022-01315-7) (2023).
- 983 31 Bishop, R. E. The bacterial lipocalins. *Biochim. Biophys. Acta* **1482**, 73-83,  
984 doi:10.1016/s0167-4838(00)00138-2 (2000).
- 985 32 Wang, J. *et al.* Gluconeogenic growth of *Vibrio cholerae* is important for competing with host  
986 gut microbiota. *J. Med. Microbiol.* **67**, 1628-1637, doi:10.1099/jmm.0.000828 (2018).
- 987 33 Ball, A. S., Chaparian, R. R. & van Kessel, J. C. Quorum Sensing Gene Regulation by  
988 LuxR/HapR Master Regulators in Vibrios. *J. Bacteriol.* **199**, doi:10.1128/jb.00105-17 (2017).
- 989 34 Merrell, D. S., Tischler, A. D., Lee, S. H. & Camilli, A. *Vibrio cholerae* requires rpoS for  
990 efficient intestinal colonization. *Infect. Immun.* **68**, 6691-6696, doi:10.1128/iai.68.12.6691-  
991 6696.2000 (2000).
- 992 35 Manneh-Roussel, J. *et al.* cAMP receptor protein controls *Vibrio cholerae* gene expression in  
993 response to host colonization. *MBio* **9**, 10.1128/mbio.00966-00918 (2018).
- 994 36 Massengo-Tiassé, R. P. & Cronan, J. E. *Vibrio cholerae* FabV defines a new class of enoyl-  
995 acyl carrier protein reductase. *J. Biol. Chem.* **283**, 1308-1316, doi:10.1074/jbc.M708171200  
996 (2008).
- 997

998 37 Slamti, L., Livny, J. & Waldor, M. K. Global gene expression and phenotypic analysis of a  
999 Vibrio cholerae rpoH deletion mutant. *J. Bacteriol.* **189**, 351-362, doi:10.1128/jb.01297-06  
1000 (2007).

1001 38 Merrell, D. S. *et al.* Host-induced epidemic spread of the cholera bacterium. *Nature* **417**, 642-  
1002 645 (2002).

1003 39 Pearcy, N. *et al.* Genome-scale metabolic models and machine Learning reveal genetic  
1004 determinants of antibiotic resistance in *Escherichia coli* and unravel the underlying metabolic  
1005 adaptation mechanisms. *mSystems* **6**, e00913-00920, doi:  
1006 <https://doi.org/10.1128/mSystems.00913-20> (2021).

1007 40 Pukatzki, S. & Provenzano, D. Vibrio cholerae as a predator: lessons from evolutionary  
1008 principles. *Front. Microbiol.* **4**, 384, doi:10.3389/fmicb.2013.00384 (2013).

1009 41 Velasco, A. M., Leguina, J. I. & Lazcano, A. Molecular evolution of the lysine biosynthetic  
1010 pathways. *J. Mol. Evol.* **55**, 445-459, doi:10.1007/s00239-002-2340-2 (2002).

1011 42 Alvarez, L., Hernandez, S. B. & Cava, F. Cell Wall Biology of Vibrio cholerae. *Annu. Rev.*  
1012 *Microbiol.* **75**, 151-174, doi:10.1146/annurev-micro-040621-122027 (2021).

1013 43 Juan, C., Torrens, G., Barceló, I. M. & Oliver, A. Interplay between Peptidoglycan Biology  
1014 and Virulence in Gram-Negative Pathogens. *Microbiol. Mol. Biol. Rev.* **82**,  
1015 doi:10.1128/mmbr.00033-18 (2018).

1016 44 Huber, M., Fröhlich, K. S., Radmer, J. & Papenfort, K. Switching fatty acid metabolism by an  
1017 RNA-controlled feed forward loop. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 8044-8054,  
1018 doi:10.1073/pnas.1920753117 (2020).

1019 45 Bekaert, M., Goffin, N., McMillan, S. & Desbois, A. P. Essential Genes of Vibrio  
1020 anguillarum and Other Vibrio spp. Guide the Development of New Drugs and Vaccines.  
1021 *Front. Microbiol.* **12**, 755801, doi:10.3389/fmicb.2021.755801 (2021).

1022 46 Brodie, A., Azaria, J. R. & Ofran, Y. How far from the SNP may the causative genes be?  
1023 *Nucleic Acids Res.* **44**, 6046-6054, doi:10.1093/nar/gkw500 (2016).

1024 47 Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**,  
1025 308-311, doi:10.1093/nar/29.1.308 (2001).

1026 48 Kastritis, P. L. & Bonvin, A. M. On the binding affinity of macromolecular interactions:  
1027 daring to ask why proteins interact. *J R Soc Interface* **10**, 20120835,  
1028 doi:10.1098/rsif.2012.0835 (2013).

1029 49 Du, X. *et al.* Insights into Protein-Ligand Interactions: Mechanisms, Models, and Methods.  
1030 *Int. J. Mol. Sci.* **17**, doi:10.3390/ijms17020144 (2016).

1031 50 Cotten, M. & Phan, M. V. T. Evolution of increased positive charge on the SARS-CoV-2  
1032 spike protein may be adaptation to human transmission. *iScience* **26**, 106230,  
1033 doi:10.1016/j.isci.2023.106230 (2023).

1034 51 Musil, M., Konegger, H., Hon, J., Bednar, D. & Damborsky, J. Computational design of  
1035 stable and soluble biocatalysts. *Acs Catalysis* **9**, 1033-1054 (2018).

1036 52 Zhou, H. X. & Pang, X. Electrostatic Interactions in Protein Structure, Folding, Binding, and  
1037 Condensation. *Chem. Rev.* **118**, 1691-1741, doi:10.1021/acs.chemrev.7b00305 (2018).

1038 53 Sperandio, V., Giron, J. A., Silveira, W. D. & Kaper, J. B. The OmpU outer membrane  
1039 protein, a potential adherence factor of Vibrio cholerae. *Infect. Immun.* **63**, 4433-4438 (1995).

1040 54 Kamp, H. D., Patimalla-Dipali, B., Lazinski, D. W., Wallace-Gadsden, F. & Camilli, A. Gene  
1041 fitness landscapes of Vibrio cholerae at important stages of its life cycle. *PLoS Pathog.* **9**,  
1042 e1003800 (2013).

1043 55 Fu, Y., Waldor, M. K. & Mekalanos, J. J. Tn-Seq analysis of Vibrio cholerae intestinal  
1044 colonization reveals a role for T6SS-mediated antibacterial activity in the host. *Cell Host*  
1045 *Microbe* **14**, 652-663, doi:10.1016/j.chom.2013.11.001 (2013).

1046 56 Provenzano, D., Lauriano, C. M. & Klose, K. E. Characterization of the role of the ToxR-  
1047 modulated outer membrane porins OmpU and OmpT in Vibrio cholerae virulence. *J.*  
1048 *Bacteriol.* **183**, 3652-3662 (2001).

1049 57 Mathur, J. & Waldor, M. K. The Vibrio cholerae ToxR-regulated porin OmpU confers  
1050 resistance to antimicrobial peptides. *Infect. Immun.* **72**, 3577-3583 (2004).

1051 58 Merrell, D. S., Bailey, C., Kaper, J. B. & Camilli, A. The ToxR-mediated organic acid  
1052 tolerance response of Vibrio cholerae requires OmpU. *J. Bacteriol.* **183**, 2746-2754 (2001).

1053 59 Li, H., Zhang, W. & Dong, C. Crystal structure of the outer membrane protein OmpU from  
1054 Vibrio cholerae at 2.2 Å resolution. *Acta Crystallographica Section D: Structural Biology* **74**,  
1055 21-29 (2018).

1056 60 Seed, K. D. *et al.* Evolutionary consequences of intra-patient phage predation on microbial  
1057 populations. *Elife* **3**, e03497, doi:10.7554/eLife.03497 (2014).

1058 61 Lim, A. N. W., Yen, M., Seed, K. D., Lazinski, D. W. & Camilli, A. A Tail Fiber Protein and  
1059 a Receptor-Binding Protein Mediate ICP2 Bacteriophage Interactions with Vibrio cholerae  
1060 OmpU. *J. Bacteriol.* **203**, e0014121, doi:10.1128/jb.00141-21 (2021).

1061 62 Silva, A. J. & Benitez, J. A. Vibrio cholerae Biofilms and Cholera Pathogenesis. *PLoS Negl.*  
1062 *Trop. Dis.* **10**, e0004330, doi:10.1371/journal.pntd.0004330 (2016).

1063 63 Weber, G. G. & Klose, K. E. The complexity of ToxT-dependent transcription in Vibrio  
1064 cholerae. *Indian J. Med. Res.* **133**, 201-206 (2011).

1065 64 Lowden, M. J. *et al.* Structure of Vibrio cholerae ToxT reveals a mechanism for fatty acid  
1066 regulation of virulence genes. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 2860-2865,  
1067 doi:10.1073/pnas.0915021107 (2010).

1068 65 Levade, I. *et al.* Predicting Vibrio cholerae Infection and Disease Severity Using  
1069 Metagenomics in a Prospective Cohort Study. *J. Infect. Dis.* **223**, 342-351,  
1070 doi:10.1093/infdis/jiaa358 (2021).

1071 66 Ravcheev, D. A., Gel'fand, M. S., Mironov, A. A. & Rakhmaninova, A. B. [Purine regulon of  
1072 gamma-proteobacteria: a detailed description]. *Genetika* **38**, 1203-1214 (2002).

1073 67 Lopez, C. M., Kovler, M. L. & Jelin, E. B. Case report of extreme gastric distention and  
1074 perforation with pathologic Sarcina ventriculi colonization and Rett syndrome. *Int. J. Surg.*  
1075 *Case Rep.* **73**, 210-212, doi:10.1016/j.ijscr.2020.07.025 (2020).

1076 68 Conner, J. G., Teschler, J. K., Jones, C. J. & Yildiz, F. H. Staying alive: Vibrio cholerae's  
1077 cycle of environmental survival, transmission, and dissemination. *Virulence mechanisms of*  
1078 *bacterial pathogens*, 593-633 (2016).

1079 69 Carey, D. E. & McNamara, P. J. The impact of triclosan on the spread of antibiotic resistance  
1080 in the environment. *Front. Microbiol.* **5**, 780, doi:10.3389/fmicb.2014.00780 (2014).

1081 70 Huang, Y. H., Lin, J. S., Ma, J. C. & Wang, H. H. Functional Characterization of Triclosan-  
1082 Resistant Enoyl-acyl-carrier Protein Reductase (FabV) in Pseudomonas aeruginosa. *Front.*  
1083 *Microbiol.* **7**, 1903, doi:10.3389/fmicb.2016.01903 (2016).

1084 71 Singh, S. M., Kongari, N., Cabello-Villegas, J. & Mallela, K. M. Missense mutations in  
1085 dystrophin that trigger muscular dystrophy decrease protein stability and lead to cross-beta  
1086 aggregates. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 15069-15074, doi:10.1073/pnas.1008818107  
1087 (2010).

1088 72 Wang, Z. & Moulton, J. SNPs, protein structure, and disease. *Hum. Mutat.* **17**, 263-270,  
1089 doi:10.1002/humu.22 (2001).

1090 73 Scheller, R. *et al.* Toward mechanistic models for genotype-phenotype correlations in  
1091 phenylketonuria using protein stability calculations. *Hum. Mutat.* **40**, 444-457,  
1092 doi:10.1002/humu.23707 (2019).

1093 74 Rakoczy, E. P., Kiel, C., McKeone, R., Stricher, F. & Serrano, L. Analysis of disease-linked  
1094 rhodopsin mutations based on structure, function, and protein stability calculations. *J. Mol.*  
1095 *Biol.* **405**, 584-606, doi:10.1016/j.jmb.2010.11.003 (2011).

1096 75 Wall, S. M. The renal physiology of pendrin (SLC26A4) and its role in hypertension.  
1097 *Novartis Found. Symp.* **273**, 231-239; discussion 239-243, 261-234 (2006).

1098 76 Harris, J. B., LaRocque, R. C., Qadri, F., Ryan, E. T. & Calderwood, S. B. Cholera. *Lancet*  
1099 **379**, 2466-2476, doi:10.1016/s0140-6736(12)60436-x (2012).

1100 77 Madi, N. *et al.* Phage predation, disease severity, and pathogen genetic diversity in cholera  
1101 patients. *Science* **384**, eadj3166, doi:10.1126/science.adj3166 (2024).

1102 78 Levade, I. *et al.* Vibrio cholerae genomic diversity within and between patients. *Microbial*  
1103 *Genomics* **3**, doi:<https://doi.org/10.1099/mgen.0.000142> (2017).

1104 79 Lypaczewski, P. *et al.* Diversity of Vibrio cholerae O1 through the human gastrointestinal  
1105 tract during cholera. *bioRxiv*, doi:10.1101/2024.02.08.579476 (2024).

1106 80 Das, R. *et al.* Vibrio cholerae in rural and urban Bangladesh, findings from hospital-based  
1107 surveillance, 2000-2021. *Sci. Rep.* **13**, 6411, doi:10.1038/s41598-023-33576-3 (2023).

1108 81 Chowdhury, F. *et al.* Concomitant Enterotoxigenic *Escherichia coli* Infection Induces  
1109 Increased Immune Responses to *Vibrio cholerae* O1 Antigens in Patients with Cholera  
1110 in Bangladesh. *Infect. Immun.* **78**, 2117-2124, doi:10.1128/iai.01426-09 (2010).

1111 82 Khan, A. I. *et al.* Epidemiology of cholera in Bangladesh: findings from nationwide hospital-  
1112 based surveillance, 2014–2018. *Clin. Infect. Dis.* **71**, 1635-1642 (2020).

1113 83 Bwire, G. *et al.* Alkaline peptone water enrichment with a dipstick test to quickly detect and  
1114 monitor cholera outbreaks. *BMC Infect. Dis.* **17**, 726, doi:10.1186/s12879-017-2824-8 (2017).

1115 84 Rahman, M. R. *et al.* A Rapid Assessment of Health Literacy and Health Status of Rohingya  
1116 Refugees Living in Cox's Bazar, Bangladesh Following the August 2017 Exodus from  
1117 Myanmar: A Cross-Sectional Study. *Tropical medicine and infectious disease* **5**, 110,  
1118 doi:10.3390/tropicalmed5030110 (2020).

1119 85 Clinical and Laboratory Standards Institute. M100 Performance standards for antimicrobial  
1120 susceptibility testing. *Clinical and Laboratory Standards Institute*, (2018).

1121 86 Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to  
1122 single-cell sequencing. *J. Comput. Biol.* **19**, 455-477 (2012).

1123 87 Lee, I. *et al.* ContEst16S: an algorithm that identifies contaminated prokaryotic genomes  
1124 using 16S RNA gene sequences. *Int. J. Syst. Evol. Microbiol.* **67**, 2053-2057 (2017).

1125 88 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-2069,  
1126 doi:<https://doi.org/10.1093/bioinformatics/btu153> (2014).

1127 89 Alcock, B. P. *et al.* CARD 2020: antibiotic resistome surveillance with the comprehensive  
1128 antibiotic resistance database. *Nucleic Acids Res.* **48**, D517-d525,  
1129 doi:<https://doi.org/10.1093/nar/gkz935> (2020).

1130 90 Seeman T. ABRicate: Mass screening of contigs for antimicrobial resistance or virulence  
1131 genes. *Github* - <https://github.com/tseemann/abricate>,  
1132 doi:<https://github.com/tseemann/abricate> (2020).

1133 91 Liu, B., Zheng, D., Jin, Q., Chen, L. & Yang, J. VFDB 2019: a comparative pathogenomic  
1134 platform with an interactive web interface. *Nucleic Acids Res.* **47**, D687-D692,  
1135 doi:<https://doi.org/10.1093/nar/gky1080> (2019).

1136 92 Carattoli, A. *et al.* In silico detection and typing of plasmids using PlasmidFinder and plasmid  
1137 multilocus sequence typing. *Antimicrob. Agents Chemother.* **58**, 3895-3903,  
1138 doi:10.1128/aac.02412-14 (2014).

1139 93 Seeman T. MLST. doi:<https://github.com/tseemann/mlst> (2022).

1140 94 Jolley, K. A. *PubMLST database*, <<https://pubmlst.org/> The PubMLST site is hosted at the  
1141 Department of Zoology, University of Oxford, UK. The site is developed and maintained by  
1142 Keith Jolley (research group of Martin Maiden). Funding is provided by The Wellcome  
1143 Trust.> (Last accessed

1144 95 Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**,  
1145 3691-3693, doi:<https://doi.org/10.1093/bioinformatics/btv421> (2015).

1146 96 Page, A. J. *et al.* SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments.  
1147 *Microbial Genomics* **2**, doi:<https://doi.org/10.1099/mgen.0.000056> (2016).

1148 97 Thorpe, H. A., Bayliss, S. C., Sheppard, S. K. & Feil, E. J. Piggy: a rapid, large-scale pan-  
1149 genome analysis tool for intergenic regions in bacteria. *Gigascience* **7**, giy015 (2018).

1150 98 Minh, B. Q. *et al.* IQ-TREE 2: New models and efficient methods for phylogenetic inference  
1151 in the genomic era. *Mol. Biol. Evol.* **37**, 1530-1534,  
1152 doi:<https://doi.org/10.1093/molbev/msaa015> (2020).

1153 99 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermin, L. S.  
1154 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods* **14**, 587-  
1155 589, doi:10.1038/nmeth.4285 (2017).

1156 100 Letunic, I. & Bork, P. Interactive Tree of Life (iTOL) v6: recent updates to the phylogenetic  
1157 tree display and annotation tool. *Nucleic Acids Res.*, doi:10.1093/nar/gkae268 (2024).

1158 101 Seeman T. snippy. doi:<https://github.com/tseemann/snippy> (2016).

1159 102 Cock, P. J. *et al.* Biopython: freely available Python tools for computational molecular  
1160 biology and bioinformatics. *Bioinformatics* **25**, 1422-1423,  
1161 doi:10.1093/bioinformatics/btp163 (2009).



1162 103 Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant bacterial  
1163 whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15-e15 (2015).

1164 104 Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif.  
1165 *Bioinformatics* **27**, 1017-1018, doi:10.1093/bioinformatics/btr064 (2011).

1166 105 Bailey, T. L., Johnson, J., Grant, C. E. & Noble, W. S. The MEME Suite. *Nucleic Acids Res.*  
1167 **43**, W39-49, doi:10.1093/nar/gkv416 (2015).

1168 106 Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence  
1169 alignments using Clustal Omega. *Mol. Syst. Biol.* **7**, 539, doi:10.1038/msb.2011.75 (2011).

1170 107 Sullivan, M. J., Petty, N. K. & Beatson, S. A. Easyfig: a genome comparison visualizer.  
1171 *Bioinformatics* **27**, 1009-1010, doi:<https://doi.org/10.1093/bioinformatics/btr039> (2011).

1172 108 Salamov, V. S. A. & Solovyevand, A. Automatic annotation of microbial genomes and  
1173 metagenomic sequences. *Metagenomics and its applications in agriculture, biomedicine and*  
1174 *environmental studies*, 61-78 (2011).

1175 109 Ebrahim, A., Lerman, J. A., Palsson, B. O. & Hyduke, D. R. COBRAPy: constraints-based  
1176 reconstruction and analysis for python. *BMC Syst. Biol.* **7**, 1-6 (2013).

1177 110 Karp, P. D. *et al.* The BioCyc collection of microbial genomes and metabolic pathways. *Brief.*  
1178 *Bioinform.* **20**, 1085-1093, doi:10.1093/bib/bbx085 (2019).

1179 111 Hagberg, A., Swart, P. & S Chult, D. Exploring network structure, dynamics, and function  
1180 using NetworkX. (Los Alamos National Lab.(LANL), Los Alamos, NM (United States),  
1181 2008).

1182 112 Hunter, J. D. Matplotlib: A 2D graphics environment. *Computing in science & engineering* **9**,  
1183 90-95 (2007).

1184 113 Ludden, C. *et al.* One health genomic surveillance of *Escherichia coli* demonstrates distinct  
1185 lineages and mobile genetic elements in isolates from humans versus livestock. *mBio* **10**,  
1186 e02693-02618, doi:10.1128/mBio.02693-18 (2019).

1187 114 Chawla, N. V., Bowyer, K. W., Hall, L. O. & Kegelmeyer, W. P. SMOTE: synthetic minority  
1188 over-sampling technique. *J. Artif. Intell. Res.* **16**, 321-357, doi:<https://doi.org/10.1613/jair.953>  
1189 (2002).

1190 115 Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825-  
1191 2830 (2011).

1192 116 Jones, E., Oliphant, T. & Peterson, P. SciPy: Open source scientific tools for Python. (2001).

1193 117 Oliphant, T. E. Python for scientific computing. *Computing in science & engineering* **9**, 10-20  
1194 (2007).

1195 118 Wainer, J. & Cawley, G. Empirical evaluation of resampling procedures for optimising SVM  
1196 hyperparameters. *J. Mach. Learn. Res.* **18**, 1-35 (2017).

1197 119 Cawley, G. C. & Talbot, N. L. On over-fitting in model selection and subsequent selection  
1198 bias in performance evaluation. *J. Mach. Learn. Res.* **11**, 2079-2107 (2010).

1199 120 Demšar, J. Statistical comparisons of classifiers over multiple data sets. *The Journal of*  
1200 *Machine Learning Research* **7**, 1-30 (2006).

1201 121 Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular  
1202 interaction networks. *Genome Res.* **13**, 2498-2504 (2003).

1203 122 Szklarczyk, D. *et al.* The STRING database in 2023: protein-protein association networks and  
1204 functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Res.* **51**,  
1205 D638-d646, doi:10.1093/nar/gkac1000 (2023).

1206 123 Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,  
1207 583-589, doi:10.1038/s41586-021-03819-2 (2021).

1208 124 Uniprot Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids*  
1209 *Res.* **51**, D523-d531, doi:10.1093/nar/gkac1052 (2023).

1210 125 Pettersen, E. F. *et al.* UCSF Chimera—A visualization system for exploratory research and  
1211 analysis. *J. Comput. Chem.* **25**, 1605-1612, doi:<https://doi.org/10.1002/jcc.20084> (2004).

1212 126 Pettersen, E. F. *et al.* UCSF ChimeraX: Structure visualization for researchers, educators, and  
1213 developers. *Protein Sci.* **30**, 70-82, doi:10.1002/pro.3943 (2021).

1214 127 Pires, D. E. V., Ascher, D. B. & Blundell, T. L. DUET: a server for predicting effects of  
1215 mutations on protein stability using an integrated computational approach. *Nucleic Acids Res.*  
1216 **42**, W314-W319, doi:10.1093/nar/gku411 (2014).

1217 128 Rodrigues, C. H., Pires, D. E. & Ascher, D. B. DynaMut: predicting the impact of mutations  
1218 on protein conformation, flexibility and stability. *Nucleic Acids Res.* **46**, W350-W355 (2018).  
1219 129 Sim, N.-L. *et al.* SIFT web server: predicting effects of amino acid substitutions on proteins.  
1220 *Nucleic Acids Res.* **40**, W452-W457 (2012).  
1221 130 Jurrus, E. *et al.* Improvements to the APBS biomolecular solvation software suite. *Protein*  
1222 *Sci.* **27**, 112-128, doi:10.1002/pro.3280 (2018).  
1223 131 tan0101/VibrioCARE: Maciel-Guerra *et al.* 2024 Nature Communications v. v1.0.0 (Zenodo,  
1224 2024).  
1225 132 R. Hijmans & University of California Berkeley Museum of Vertebrate Zoology. *First-level*  
1226 *Administrative Divisions, Bangladesh, 2015*, <<https://purl.stanford.edu/ps807dh8348>> (2015).  
1227 Last accessed 12 Aug 2024

1228

## 1229 **Acknowledgements**

1230 This study was supported by Research England grant [H53802] as part of the Internal Global  
1231 Challenges Research Fund Award from the University of Nottingham (TD), the Turkish Ministry of  
1232 National Education (KB) and UKRI MRC grant (MR/X009246/1) (MB).

## 1233 **Author contributions**

1234 Designed and supervised the study: M.M.H., Z.H.H, T.S., F.Q. and T.D.

1235 Planned the methodology: M.M.H., Z.H.H, N.Se., T.S., and T.D.

1236 Writing—original draft: K.B, A.M.G., M.B., and T.D.

1237 Writing—review and editing: K.B, A.M.G., M.B., N.Se., and T.D.

1238 Carried out the experiments and collected the samples: A.R., M.M.H., A.S., J.A., S.U.,  
1239 M.F.R.S., N.Su., A.I.K., Y.A.B., M.H.A., Z.H.H., T.S., F.Q.

1240 Data analysis and visualization: K.B, A.M.G., and M.B.

1241 Acquired funding: M.M.H., T.S., F.Q., and T.D.

1242

## 1243 **Competing Interests**

1244

1245 The authors declare no competing interests.

## 1246 **Figure Legends**

1247

1248 **Figure 1.** Maximum likelihood phylogenetic tree of the whole cohort based on the core genome of  
1249 129 isolates cultured from in-patients admitted to hospitals in six districts (Barisal, Chittagong,  
1250 Dhaka, Khulna, Rajshahi and Sylhet) of Bangladesh. The two distinct BD-1.2 and BD-2 lineages are  
1251 shown in the inner ring. The outer rings display additional information including serotypes, year of  
1252 collection, presence of *Vibrio* pathogenicity island VPI2 variants, presence of *Vibrio* seventh  
1253 pandemic island II (VSP2) variants, presence of phage-inducible chromosomal island-like elements 1  
1254 and 2 (PLE) and region of collection. A map of Bangladesh<sup>132</sup> showing the proportion of samples  
1255 included from each regional division is also shown.

1256 **Figure 2.** SNP network analysis of highly connected isolates. Network diagram showing pairwise  
1257 connections between isolates in our cohort with less than 15 pairwise single nucleotide polymorphisms  
1258 (SNP) differences. The panels show the same network with the nodes colour-coded according to (A)  
1259 lineages, (B) year of collection, (C) serotypes and (D) location of collection. The lines between pairs of  
1260 isolates are colour-coded by single nucleotide SNP number.

1261 **Figure 3.** An overview of the metabolic pathways associated to the core genes underlying the BD-1.2  
1262 and BD-2 lineages separation. All genes annotated were found to have reduced flux span through the  
1263 metabolic system when knocked out. Genes coloured in blue have a significant different allelic  
1264 distribution between BD-1.2 and BD-2, associated metabolic pathways are labelled in purple. All 3D  
1265 protein structures were generated in Alphafold<sup>123</sup> under a Creative Commons Attribution 4.0 license  
1266 ([CC-BY 4.0](#)), no changes were made.

1267 **Figure 4.** Supervised machine learning pipeline accurately predicts the clinical manifestations of  
1268 hospitalized patients from the genomic determinants extracted from BD-1.2 isolates, collected among  
1269 the same hospitalised patients. (A) Flow diagram showing machine learning pipeline including data  
1270 (green), pre-processing steps (yellow) and classification (blue). (B) Machine learning performance  
1271 results measured by the area under the curve (AUC) from 30 training runs for clinical symptom  
1272 combination. The results shown are for the best classifier Logistic Regression, as defined by the  
1273 Nemenyi test (Fig. S18). The violin plots show the distribution of the data, with each data point  
1274 representing one classification model. Inside each violin plot is a box plot, with the box showing the

1275 interquartile range (IQR), the whiskers showing the rest of the distribution as a proportion of 1.5 x IQR  
1276 and the white circle representing the median value. (C) Number of features (accessory genes, core  
1277 genome SNPs and intergenic SNPs) selected for each symptom. Predictive models were generated for  
1278 six different clinical symptoms (X axis): abdominal pain; dehydration Moderate vs Severe; duration of  
1279 diarrhoea <1 day vs. 1-3 days; number of stools 11-15 times vs. 16-20 times; number of stools 11-15  
1280 times vs. 21+ times; and vomit.

1281 **Figure 5.** Undirected graph network illustrating the genomic features associated with clinical symptom  
1282 models for *V. cholerae*. Node colour denotes the genomic determinant category, (i.e. accessory genes  
1283 and/or core genome coding, and intergenic SNPs) identified by machine learning. Nodes are labelled  
1284 with numbers corresponding to specific genes associated with each genomic determinant, as detailed in  
1285 the Genes legend, while unnumbered nodes are related to unannotated (hypothetical) genes. The clinical  
1286 symptom models are highlighted in different colours and explained in the legend Symptoms legend  
1287 featuring abdominal pain; dehydration Moderate vs Severe; duration of diarrhoea <1 day vs. 1-3 days;  
1288 number of stools 11-15 times vs. 16-20 times; number of stools 11-15 times vs. 21+ times; and vomit.

1289 **Figure 6.** An overview of the metabolic pathways impacted by statistically significant genes underlying  
1290 clinical symptoms. All genes annotated were found to have reduced the flux span through the metabolic  
1291 system when knocked out. Genes coloured in pink and purple carried mutations or are accessory genes  
1292 associated to the clinical symptom, respectively, and connected metabolic pathways (labelled in blue).  
1293 The genes coloured in purple were also found as statistically significant in differentiating the BD-2 and  
1294 BD-1.2 lineages (see previous sections). All 3D protein structures were generated in AlphaFold<sup>123</sup> under  
1295 a Creative Commons Attribution 4.0 license ([CC-BY 4.0](https://creativecommons.org/licenses/by/4.0/)), no changes were made.

1296 **Figure 7.** 3D protein structure analysis of FabV allelic variants underlying BD-1.2 and BD-2 lineage  
1297 evolution and clinical symptoms. (A) Violin plot indicating the distribution of the diarrhoea duration  
1298 score (0: no diarrhoea, 1: <1day, 2: 1-3 days, 3: 4-6 days and 4: 7-9 days) for the isolates containing  
1299 either Pro149 (P) or His149 (H). Statistical significance was tested with a two-sided Mann Whitney U  
1300 test, p-value is shown. (B) Violin plot indicating the distribution of the number of stools score (0: <3

1301 times, 1: 3-5 times; 2: 6-10 times; 3: 11-15 times; 4: 16-20 times; 5: 21+ times) for the isolates  
1302 containing either Pro149 (P) or His149 (H). Statistical significance was tested with a two-sided Mann  
1303 Whitney U test, p-value is shown. (C) The bar graph displays the number of isolates in the two BD  
1304 lineages associated with Pro149 (P) and His149 (H). (D) 3D structures of FabV (AlphaFold) with  
1305 Pro149 and coloured by functional domains. Amino acid residues (Lys148, Ser151, and Trp159)  
1306 interacting with Pro149 (green) are shown in sticks models. (E) 3D structures of FabV (AlphaFold)  
1307 with His149 and coloured by functional domains. Amino acid residues (Lys148, Arg 150, Ser151, and  
1308 Trp159) interacting with His149 (orange) are shown in sticks models.

1309 **Figure 8.** 3D protein structure analysis of GshB allelic variants underlying BD-1.2 and BD-2 lineage  
1310 evolution and clinical symptoms. (A) Violin plot indicating the distribution of the diarrhoea duration  
1311 score (0: no diarrhoea, 1: <1day, 2: 1-3 days, 3: 4-6 days and 4: 7-9 days) for the isolates containing  
1312 either Thr93 (T) or Ile93 (I). Statistical significance was tested with a two-sided Mann Whitney U test,  
1313 p-value is shown. (B) Violin plot indicating the distribution of the number of stools score (0: <3 times,  
1314 1: 3-5 times; 2: 6-10 times; 3: 11-15 times; 4: 16-20 times; 5: 21+ times) for the isolates containing  
1315 either Thr93 (T) or Ile93 (I). Statistical significance was tested with a two-sided Mann Whitney U test,  
1316 p-value is shown. (C) The bar graph displays the number of isolates in the two BD lineages associated  
1317 Thr93 (T) or Ile93 (I) (D) 3D structures of GshB (AlphaFold) with Thr93 and coloured by functional  
1318 domains. Amino acid residues (Asp92, Ile96, and Tyr97) interacting with Thr93 (green) are shown in  
1319 sticks models. (E) 3D structures of GshB (AlphaFold) with Ile93 and coloured by functional domains.  
1320 Amino acid residues interacting with Ile93 (orange) are shown in sticks models.

1321