

Long-Read-Sequenced Reference Genomes of the Seven Major Lineages of Enterotoxigenic Escherichia Coli (ETEC) Circulating in Modern Time

Astrid von Mentzer (✉ avm@sanger.ac.uk)

Wellcome Sanger Institute

Grace A Blackwell

European Bioinformatics Institute

Derek Pickard

University of Cambridge

Christine J Boinett

Wellcome Sanger Institute

Enrique Joffré

Karolinska Institute

Andrew J Page

Quadram Institute

Ann-Mari Svennerholm

University of Gothenburg

Gordon Dougan

University of Cambridge

Åsa Sjöling

Karolinska Institute

Research Article

Keywords: Reference genome, plasmid, enterotoxigenic Escherichia coli, ETEC, diarrhoea

Posted Date: February 25th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-237525/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Scientific Reports on April 29th, 2021. See the published version at <https://doi.org/10.1038/s41598-021-88316-2>.

1 **Long-read-sequenced reference genomes of the seven major lineages**
2 **of enterotoxigenic *Escherichia coli* (ETEC) circulating in modern time**

3

4 **Astrid von Mentzer^{1,2*#}, Grace A. Blackwell^{1,3}, Derek Pickard⁴, Christine J.**
5 **Boinett¹, Enrique Joffré⁵, Andrew J Page^{1,6}, Ann-Mari Svennerholm², Gordon**
6 **Dougan⁴ and Åsa Sjöling⁵**

7

8 ¹Wellcome Sanger Institute, Hinxton, Cambridge, UK

9 ²Department of Microbiology and Immunology, Sahlgrenska Academy, University of
10 Gothenburg, Sweden

11 ³EMBL-EBI, Hinxton, Cambridge, UK

12 ⁴University of Cambridge, Department of Medicine, Cambridge, UK

13 ⁵ Department of Microbiology, Tumor and Cell Biology, Karolinska Institutet

14 ⁶Quadram Institute Bioscience, Norwich Research Park, Norwich, UK

15 *corresponding author: avm@sanger.ac.uk or mentzerv@chalmers.se

16 #Currently affiliated with Chalmers University of Technology, Gothenburg, Sweden

17

18 **Abstract**

19 Enterotoxigenic *Escherichia coli* (ETEC) is an enteric pathogen responsible for the
20 majority of diarrheal cases worldwide. ETEC infections are estimated to cause 80,000
21 fatalities per year, with the highest rates of burden, ca 75 million cases per year,
22 amongst children under five years of age in resource-poor countries. It is also the
23 leading cause of diarrhoea in travellers. Previous large-scale sequencing studies have
24 found seven major ETEC lineages currently in circulation worldwide.

25

26 We used PacBio long-read sequencing combined with Illumina sequencing to create
27 high-quality complete reference genomes for each of the major lineages with manually
28 curated chromosomes and plasmids. We confirm that the major ETEC lineages all
29 harbour conserved plasmids that have been associated with their respective background
30 genomes for decades and that the plasmids and chromosomes of ETEC are both crucial
31 for ETEC virulence and success as pathogens. The in-depth analysis of gene content,
32 synteny and correct annotations of plasmids will elucidate other plasmids with and
33 without virulence factors in related bacterial species. These reference genomes allow
34 for fast and accurate comparison between different ETEC strains, and these data will
35 form the foundation of ETEC genomics research for years to come.

36

37 **Keywords:** Reference genome, plasmid, enterotoxigenic *Escherichia coli*, ETEC,
38 diarrhoea

39

40 **Introduction**

41 Diarrheal pathogens are a leading cause of morbidity and mortality globally (WHO,
42 2017), with enterotoxigenic *Escherichia coli* (ETEC) accounting for a large proportion
43 of the diarrhoea cases in resource-poor countries [1]. An estimation of 220 million cases
44 each year are attributed to ETEC (WHO PPC 2020). The most vulnerable group is
45 children under five years, but ETEC can also cause disease in adults and is the principal
46 cause of diarrhoea in travellers. Resource-poor settings, where access to clean water is
47 limited, enable the spread of ETEC, transmitted via the faecal-oral route through
48 ingestion of contaminated food or water [2]. The disease severity may range from mild
49 to cholera-like symptoms with profuse watery diarrhoea. The infection is usually self-
50 limiting, lasting three to four days and may be treated by water and electrolyte

51 rehydration to balance the loss of fluids and ions. There is strong evidence to support
52 that an ETEC vaccine is of key importance to prevent children and adults from
53 developing ETEC disease [3]. Several efforts are on-going to develop an ETEC
54 vaccine, with the majority focusing on including immunogenic antigens possibly
55 capable of inducing protection against a majority of the circulating ETEC clones [4–6].

56

57 ETEC bacteria adhere to the small intestine through fimbrial, fibrillar or afimbrial outer
58 membrane-structures called colonisation factors (CF). Upon colonisation, the bacteria
59 proliferate and secrete heat-labile toxin (LT) and/or heat-stable toxins, (STh or STp)
60 causing diarrhoea and often vomiting causing the further spread of the bacteria in the
61 environment [7].

62

63 The ability of an ETEC strain to infect relies on its ability to adhere to cells of a specific
64 host. To this day, 27 different CFs with human tropism have been described, and
65 individual ETEC strains usually express 1-3 different CFs [8–14]. The enterotoxins, LT
66 and ST, can also be subdivided based on structure and function. Human-associated
67 ETEC strains express one of the 28 different LT-I variants (LTh-1 and LTh-2 are the
68 most common variants) [15] alone or together with one of the genetically distinct types
69 of STa; STh and STp [16,17].

70

71 We have previously shown that ETEC strains causing human disease can be grouped
72 into a set of clonal lineages that encompass strains with specific virulence profiles.
73 Seven of the 21 identified lineages encompass ETEC strains that express the most
74 commonly found CFs and toxin profiles amongst isolated clinical ETEC strains [4,18].

75

76 There is currently one complete ETEC reference genome, H10407 [19], with curated
77 annotations. Several additional complete ETEC genomes are available [20,21], some
78 of which are annotated using automated annotation pipelines that often fail at correctly
79 annotating ETEC specific genes such as CFs. The rapid adaptation of next-generation
80 sequencing in public health, specifically within bacterial diseases [22,23] and several
81 large-scale sequencing studies [18,24–28] has led to a sharp increase in the number of
82 publicly available ETEC genomes. Most of these data were generated with short-read
83 technologies, such as Illumina. A limitation of short-read sequence data is the inability
84 to unambiguously resolve repetitive regions of a genome, leading to fragmented de
85 novo assemblies of the underlying genome, missing regions and genes, and disjointed
86 synteny. ETEC is a highly diverse pathogen both in the core genome (SNPs) and the
87 accessory genome, including mobile genetic elements (MGE). Clinically related
88 MGEs, such as virulence plasmids, vary within ETEC strains. Hence, it is important to
89 identify lineage-specific reference genomes that are carefully annotated, i.e. manually
90 curated annotations, and include both chromosome and plasmid(s). Several complete
91 genomes have been generated using long-read sequencing alone [20,28], however,
92 circularising some chromosomes and plasmids may be difficult, and small plasmids can
93 be lost. Assembly issues can be resolved using a hybrid assembly approach combining
94 long-read and short-read sequencing data. In this report, we describe eight genomes,
95 eight chromosomes (seven successfully circularised) and 29 plasmids (24 successfully
96 circularised) with curated annotations, from isolates representing the major ETEC
97 lineages (L1-L7) that cause disease globally. They are sequenced using both short and
98 long-read sequencing technologies to provide the highest accuracy currently available.
99 These reference genomes will form the foundation of ETEC genomics research for
100 years to come.

101

102 **Results**

103 ***Genome analysis of eight representative ETEC isolates***

104 Eight ETEC strains representing the seven major ETEC lineages (L1-L7) comprising
105 isolates with the most prevalent virulence factor profiles were sequenced, assembled,
106 circularised and manually curated (Table 1).

107 L3 includes two different representative strains, one CS7 and one CFA/I positive strain.
108 All chromosomes except one were circularised (E1779). The average length of the
109 chromosome was 4,927,521 bases (4,721,269-5,151,162) with an average GC content
110 of 50.7% (50.4% to 50.9%) and the number of CDS ranging from 4,409 to 4,924 (Table
111 S1). Each ETEC reference genome contains between two and five plasmids
112 encompassing plasmid-specific features. Some of which carried virulence genes and/or
113 antibiotic resistance genes (Table 2, Additional File 2).

114

115 ***Comparative genomics of the chromosome***

116 The chromosomes of the reference strains were aligned and compared using MAUVE,
117 and the overall structure is conserved across all eight chromosomes (Figure S1). In
118 total, 8,348 chromosomal genes were identified in the eight ETEC strains with 3,179
119 genes considered part of the core genome shared by all eight reference strains. The
120 majority of human commensal *Escherichia coli* (*E. coli*) strains belong to subgroup A
121 [29,30]. However, ETEC strains fall into multiple phylogenetic groups (A, B1, B2, D,
122 E, F and CladeI with the majority found in the phylogenetic groups A and B1 [18]. The
123 phylogenetic group of the eight ETEC reference strains have previously been
124 determined using the triplex-PCR scheme [31]. The ETEC references were re-analysed
125 using ClermonTyping [32] and it was determined that strain E1373 belongs to the

126 phylogenetic group E while the other reference isolates belong to groups A and B1
127 (Table 1).

128

129 ***Plasmids***

130 The plasmids of each isolate were annotated using Prokka followed by manual curation
131 of the annotations including genes part of the conjugation machinery and known
132 plasmid stability genes. Virulence factors (including CFs, toxins and putative virulence
133 factors), antibiotic resistance determinants with the Comprehensive Antibiotic
134 Resistance Database (CARD) [33] as well as complete and partial insertion elements
135 and prophages were manually annotated. The plasmids were designated
136 pAvM_strainID_integer, e.g. pAvM_E925_4 (Additional file 2). The first plasmid
137 reported in this study starts at 4 as three previous plasmids E873p1-3 already have been
138 deposited to GenBank related to a different project [8].

139

140 Plasmids were typed by analysing the presence and variation of specific replication
141 genes to assign the plasmids to incompatibility (Inc) groups. The Inc groups of the
142 ETEC reference plasmids were first determined using PlasmidFinder and further
143 classified into subtypes using pMLST [34]. The replicons identified are IncFII,
144 IncFIIA, IncFIIS, IncFIB, IncFIC, IncI1 and IncY. Plasmids with replicon IncY,
145 IncFIIY or IncB/O/K/Z mainly harboured plasmid associated genes, such as stability
146 and transfer genes. Importantly, replicons FII, FIB and I1 were found to be strongly
147 associated with virulence genes as all CF loci, toxins and virulence factors *eatA* and
148 *etpBAC* were present on these plasmids. The majority of all ETEC plasmids analysed
149 here (17/29) belong to IncFII, of which six of the IncFII plasmids have an additional
150 IncFIB replicon. In six of the ETEC reference strains two or three IncFII replicons are

151 present, for example, in strain E925, the plasmids pAvM_E925_4 and 7 both belong to
152 IncFII. However, the plasmids were further subtyped to FII-111 and FII-15,
153 respectively, (Table 2 and Additional file 3), explaining the plasmid compatibility.

154

155 ***Virulence factors***

156 The CFs expressed by the selected reference strains are CFA/I, CS1-CS3, CS5-CS7 and
157 CS21. Three of the strains (E925, E1649 and E1779) express both LT and ST, two
158 strains (E2980 and E1441) express LT and the strains E36 and E562 express STh, while
159 E1373 express STp (Table 1). A plasmid can harbour multiple virulence genes, usually
160 a CF locus and genes encoding one or two toxins. Interestingly, plasmids do not often
161 harbour multiple CF loci, but on individual plasmids (in the ETEC reference strains
162 described here). Exceptions for this is strain E1779 in which CS5 and CS6 loci are
163 located on the same plasmid (pAvM_E1779_19). In both E925 (L1) and E1649 (L2)
164 the genes encoding CS3 (*cstABGH*), ST (*estA*) and LT (*eltAB*) are located on the same
165 plasmid, both with the FII replicon and of roughly the same size (Table 2). Blastn
166 comparison between the plasmids and additional plasmids that harbour the same
167 virulence genes shows that they are highly conserved (Figure 1). The results correspond
168 with the close genetic relationship and common ancestry of lineage 1 (L1) and lineage
169 2 (L2) [18].

170

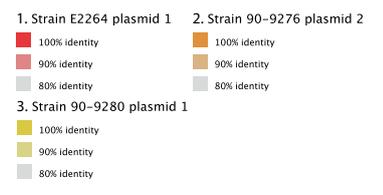
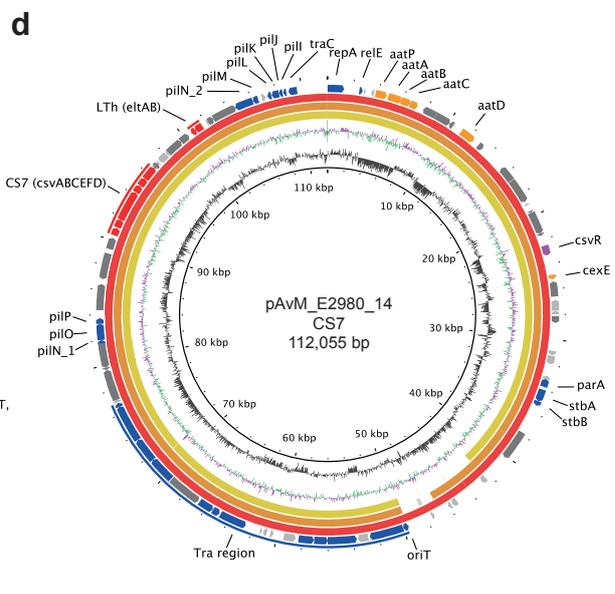
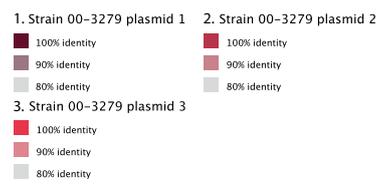
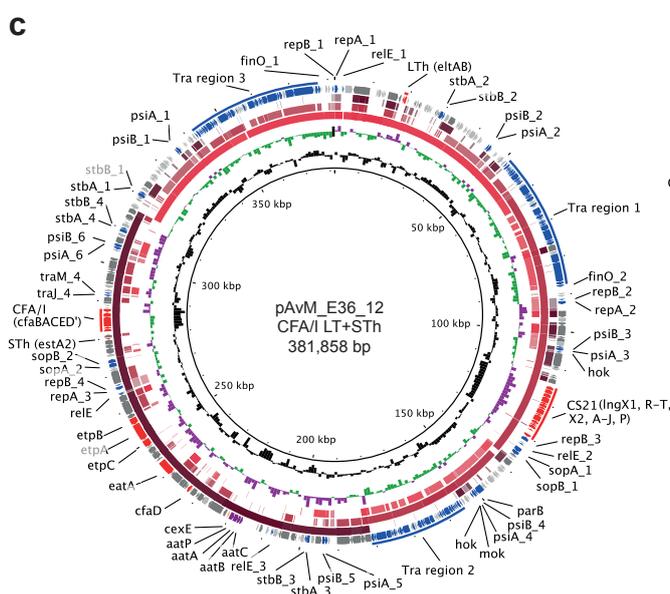
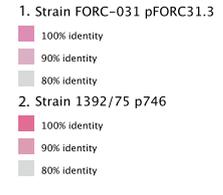
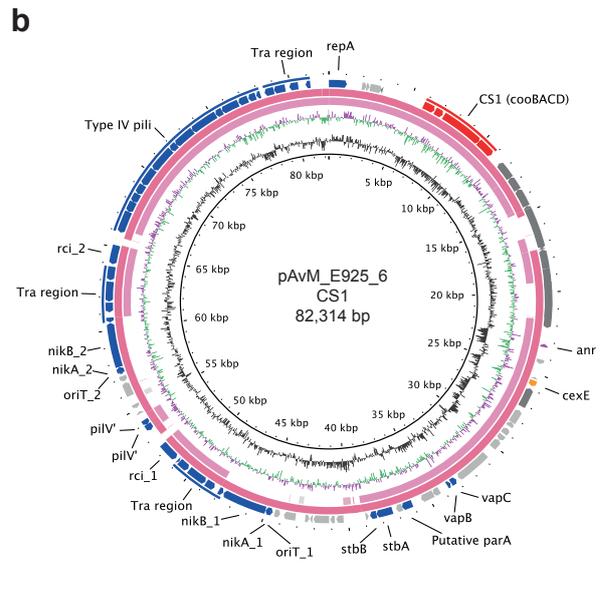
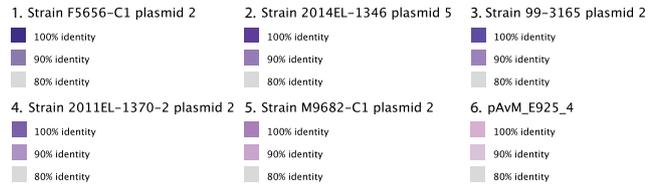
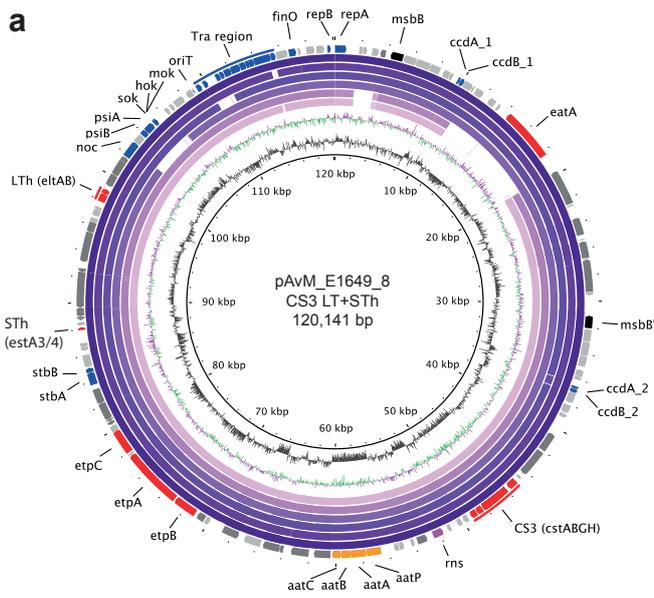
171

172

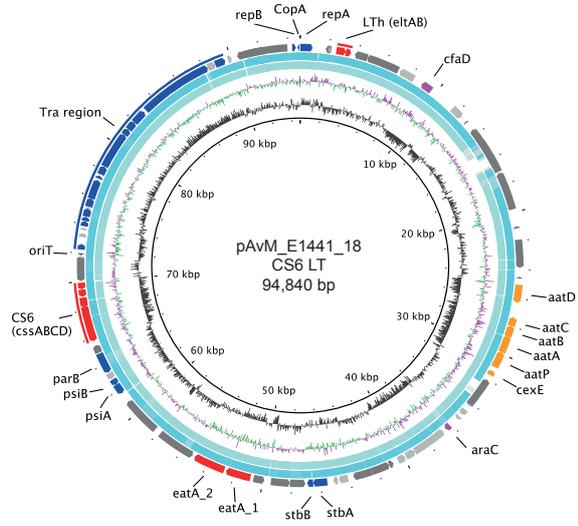
173

174

175

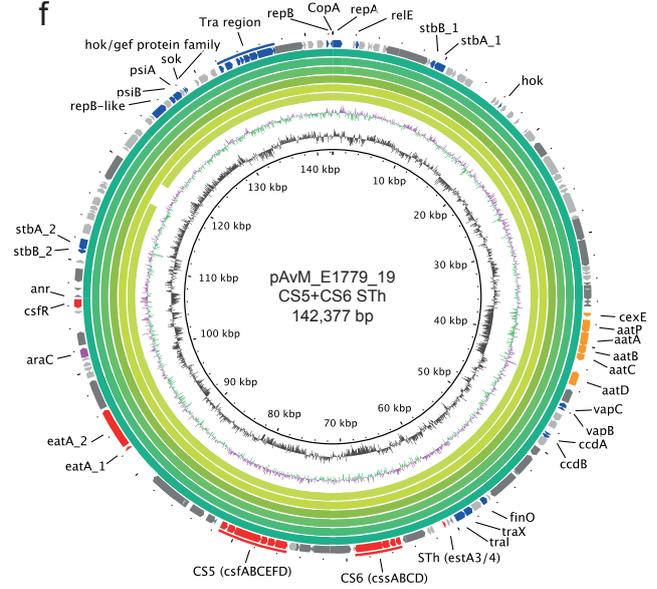


e



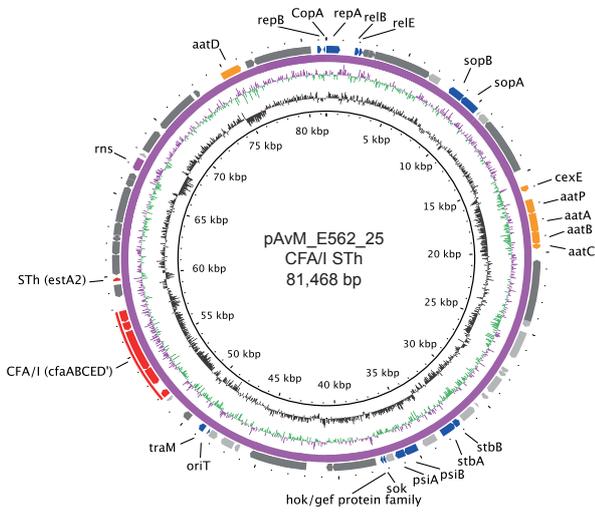
1. Strain F5505-C1 plasmid 2
 - 100% identity
 - 90% identity
 - 80% identity
2. Strain ATCC 43886 plasmid 2
 - 100% identity
 - 90% identity
 - 80% identity

f



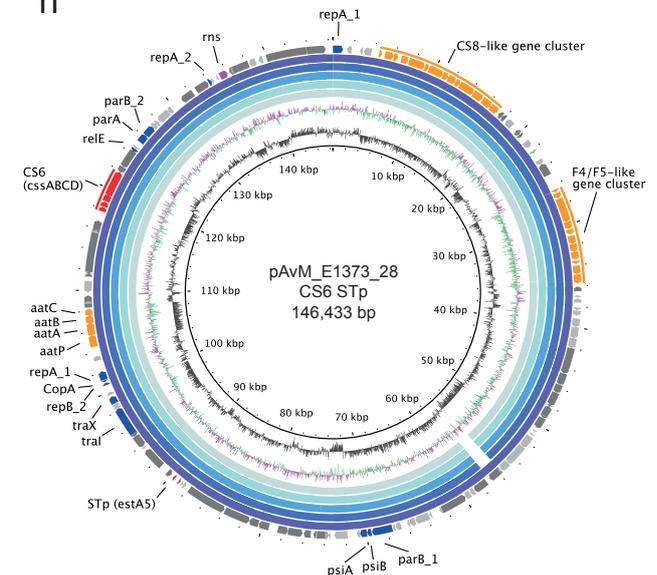
1. Strain 204576 p146
 - 100% identity
 - 90% identity
 - 80% identity
2. Strain 120899 p146
 - 100% identity
 - 90% identity
 - 80% identity
3. Strain E2265 plasmid 1
 - 100% identity
 - 90% identity
 - 80% identity
4. Strain 504237 p142
 - 100% identity
 - 90% identity
 - 80% identity
5. Strain 602354 p148
 - 100% identity
 - 90% identity
 - 80% identity
6. Strain F5176-C6 plasmid 1
 - 100% identity
 - 90% identity
 - 80% identity

g

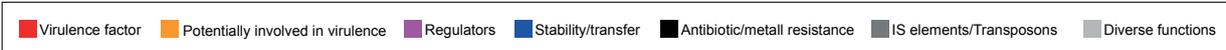


1. Strain 504239 p101
 - 100% identity
 - 90% identity
 - 80% identity

h



1. Strain F8111-1SC3 plasmid 3
 - 100% identity
 - 90% identity
 - 80% identity
2. pEntYN10
 - 100% identity
 - 90% identity
 - 80% identity
3. Strain F9792 plasmid
 - 100% identity
 - 90% identity
 - 80% identity
4. Strain 2014EL-1345-2 plasmid 4
 - 100% identity
 - 90% identity
 - 80% identity
5. Strain F6326-C1 plasmid 2
 - 100% identity
 - 90% identity
 - 80% identity



176 **Figure 1:** Comparison between the ETEC reference plasmids harbouring colonisation
177 factors and other PacBio-sequenced ETEC plasmids using blastn.

178 **a)** pAvM_E1649_8 (CS3) as reference and pAvM_925_4 (CS3) compared to the
179 following ETEC plasmids: F5656-C1 plasmid 2 (USA, CP024262.1), 2014-EL-1346-
180 6 plasmid 5 (2014, USA; CP024237.1), 99-3165 plasmid 2 (USA; CP029980.1),
181 2011EL-1370-2 plasmid 2 (2011, USA; CP022914.1) and M9682-C1 plasmid 2 (1975,
182 USA; CP024277.1). **b)** pAvM_E925_6 (CS1) compared to ETEC plasmids
183 pFORC31.3 (2004, Korea; CP013193.1) and 1392/75 p746 (1973; FN822748.1). **c)**
184 pAvM_E36_12 (CFA/I) compared to plasmids 1-3 (p1: CP024294.1; p2 CP024295.1;
185 p3: CP024296.1) from ETEC strain 00-3279 (USA). **d)** pAvM_E2980_14 (CS7)
186 compared to E2264 plasmid 1 (2006, Bangladesh; CP023350.1), 90-9276 plasmid 2
187 (1988, Bangladesh; CP024298.1) and 90-9280 plasmid 1 (1988, Bangladesh;
188 CP024241.1). **e)** pAvM_E1441_18 (CS6) compared to F5505-C1 plasmid 2 (2013,
189 Sweden; CP023259.1) and ATCC 43886 plasmid 2 (CP024255.1). **f)**
190 pAvM_E1779_19 (CS5+CS6) compared to 204576 p146 (2010, Mali; CP025908.1),
191 120899 p146 (2012, Gambia; CP025917.1), E2265 plasmid 1 (2006, Bangladesh;
192 CP023347.1), 504237 p142 (2010, India; CP025863.1), 602354 p148 (2009,
193 Bangladesh; CP025848.1) and F5176-C6 plasmid 1 (1997; CP024668.1). **g)**
194 pAvM_E562_25 (CFA/I) compared to p504239_101 (2010, India; CP025860.1). **h)**
195 pAvM_E1373_28 (CS6) compared to F8111-1SC3 plasmid 3 (USA; CP024272.1),
196 pEntYN10 (1991, Japan; AP014654.2), F9792 plasmid (USA; CP023274.1), 2014EL-
197 1345-2 plasmid 4 (2014, USA; CP024227.1) and F6326-C1 plasmid 2 (1998, USA;
198 CP024265.1). The thresholds chosen for the blastn are shown in the key below each
199 plasmid comparison. The colour code for the annotations are listed at the bottom of the
200 figure. The two most inner rings depict GC content in black and GC Skew- in purple
201 and GC Skew+ in green. The figures were generated using BRIG [84] v0.95.
202

203 Besides CFs and toxins, additional virulence factors were identified in the majority of
204 the strains (Table 2), with *eatA* and *etpBAC* being the most commonly found.

205 EatA is an immunogenic mucinase that contributes to virulence by degrading MUC2
206 which is the major protein component of mucus in the small intestine [35,36]. The
207 *etpABC* genes encode an adhesin located on the tip of the flagella and mediate
208 adherence to host cells [37,38]. Four reference strains (E925, E1649, E36 and E562)
209 harbour both *eatA* and *etpBAC*. In three strains the *eatA* and/or *etpBAC* are located on
210 the same plasmid with an FII or FII+FIB replicon along with additional ETEC virulence
211 genes, except in E562 and E1373, where *eatA* and *etpBAC* are located on an I1+FII

212 (pAvM_E562_23) and I1 (pAvM_E1373_16) plasmid, respectively, which mainly
213 contains plasmid associated genes including genes encoding the *pil* operon and *tra*-
214 operon (pAvM_E562_23). Furthermore, a less explored putative virulence factor is
215 CexE, which is an extracytoplasmic protein dependent on the expression of the CFA/I
216 regulator *cfad* [39], and was first identified in H10407 [40]. Corroborating earlier
217 findings, the CFA/I positive E36 (L3) and E562 (L6) isolates harbour *cexE*
218 (pAvM_E36_12 and pAvM_E562_25). In addition, *cexE* is present in pAvM_E925_6,
219 pAvM_E1779_19 and pAVM_E2980_14, pAvM_E1441_18 and pAvM_E1373_28.
220 *CexE* has previously also been identified in several CS5+CS6 positive ETEC and
221 shown to be upregulated in the presence of bile and sodium glycocholate-hydrate [41].
222 Bile is known to be involved in the regulation of several ETEC CFs [42,43]. The
223 location of *cexE* seems to be conserved across specific strains. In pAvM_E36_12,
224 pAvM_E1441_18, pAvM_E1779_19 and pAvM_E562_25 *cexE* is located upstream of
225 the *aatPABC* locus, whereas in pAvM_E925_6 and pAvM_E2980_14 *cexE* is located
226 downstream of *rob* (an AraC family transcriptional regulator) in the opposite direction.
227 pAvM_E925_4 harbours the *aatPABC* locus; however *cexE* is located on a different
228 plasmid (pAvM_E925_6) in this strain.

229

230 ***Comparison of plasmids with the same virulence profile***

231 ETEC isolates within a lineage share the same virulence profile, specifically the same
232 CF profile (Figure S10). We verified that our selected isolates grouped within
233 previously described lineages with confirmed virulence profiles by phylogenetic
234 analyses (Figure S2). Blastn of each of the CF positive plasmids from each reference
235 genome were performed, and the best hit(s) were used for subsequent analysis (Figure
236 1). Most of the plasmids identified as related to the ETEC reference plasmids were not

237 annotated, hence, when needed these were annotated using the corresponding ETEC
238 reference plasmids annotation as a high priority when running Prokka. We show that
239 plasmids with the same CF and toxin profile from the same lineage are often conserved
240 (Figure 1). For example, the two plasmids encoding CS3 (pAvM_E925_4 and
241 pAvM_E1649_8) are highly similar to several CS3 harbouring plasmids from O6:H16
242 strains collected from various geographical locations between 1975 and 2014, including
243 *E. coli* O6:H16 strain M9682-C1 plasmid unnamed2 (CP024277.1) and *E. coli* strain
244 O6:H16 F5656C1 plasmid unnamed2 (CP024262.1) PacBio sequenced by Smith et al.
245 [20] (Figure 1a). Furthermore, high coverage and similarity were found between the
246 plasmids of isolates E1441 (L4), and PacBio sequenced plasmids of ETEC isolates
247 ATCC 43886/E2539C1 and 2014EL-1346-6 [20]. These isolates were collected in the
248 seventies [44] and 2014 (from a CDC collection), respectively, and assigned as O25:
249 H16 which is the O group determined for E1441 in silico (Figure 1e). Plasmids of
250 E2980 (LT+ CS7, L3) were validated by the PacBio sequenced plasmids of ETEC
251 isolate E2264 (Figure 1d). Similarly, two plasmids of E1779 (LT, STh+CS5+CS6, L5)
252 was identified in E2265 (LT, STh+CS5+CS6 [28,41], although E1779 harboured two
253 additional plasmids. Several additional L5 ETEC genomes have been sequenced within
254 the GEMS study [45], and high plasmid similarity and conservation in CS5+CS6
255 positive L5 isolates was evident (Figure 1f).

256

257 Overall the results show that ETEC plasmids are specific to lineages circulating
258 worldwide and conserved over time (Figure 1, Figure S2-S9 for more extensive plasmid
259 annotation and Figure S10). Thus, the plasmids of major ETEC lineages must confer
260 evolutionary advantages to their host genomes since they are seldom lost.

261

262 ***Antibiotic resistance***

263 *E. coli* can become resistant to antibiotics, both via the presence of antibiotic resistance
264 genes and the acquisition of adaptive and mutational changes in genes encoding efflux
265 pumps and porins which allows the bacterium to pump out the antibiotic molecules
266 effectively.

267

268 Antibiotic resistance genomic marker(s), both chromosomally located and on plasmids,
269 were identified using the CARD database [33] (Table 2, Figure S12 and Figure S13 and
270 Additional file 2). Similar to other studies, IncFII and B/O/K/Z plasmids were found to
271 harbour genes conferring antibiotic resistance [46]. Furthermore, the phenotypic
272 antibiotic resistance profile was determined with clinical MIC breakpoints based on
273 EUCAST (The European Committee on Antimicrobial Susceptibility Testing) [47]
274 (Table S2). Phenotypic antibiotic resistance profiles (Table S2) were supported mainly
275 by the findings of antibiotic resistance genes, efflux pumps and porins (Figures S4 and
276 S5 and Table S3), although some differences were found. All ETEC reference strains
277 are phenotypically resistant to at least two antibiotics of the 14 tested (Table S2).
278 Resistance against penicillin's, norfloxacin (Nor) and chloramphenicol (Cm) is most
279 common among these strains. Two of the strains, E1441 and E2980, harbour more than
280 four antibiotic resistance genes as well as multiple efflux systems and porins (Figure
281 S12, Figure S13 and Table S3). The plasmid pAvM_E1441_17 carries *aadA1-like*,
282 *dfrA15*, *sul1* and *tetA(A)* resistance genes (Table 2), where the first three genes are in a
283 Class 1 integron which confers resistance to streptomycin, trimethoprim, and
284 sulphonamide (sulphamethoxazole). The gene *tetA(A)* is part of a truncated Tn1721
285 transposon [48]. The E1441 strain was verified as resistant to tetracycline (Tet) and
286 sulphamethoxazole-trimethoprim (Sxt) while streptomycin was not tested. A *mer*

287 operon derived from Tn21 is also present the resistance region of pAvM_E1441_17
288 (Table 2), indicating that the plasmid would also likely confer tolerance to mercury,
289 although this was not confirmed. Interestingly, this multi-replicon (FII and FIB)
290 plasmid also harbours the *lng* locus encoding CS21, one of the most prevalent ETEC
291 CFs. In isolate E2980 virulence plasmid pAvM_E2980_15 harboured multiple
292 resistance genes in the same region (*bla*_{TEM-1b}, *strA*, *strB* and *sul2*) conferring resistance
293 to ampicillin, streptomycin and sulphonamides. E2980 was found to be resistant to
294 ampicillin (Amp) and oxacillin (Oxa), which can be broken down by the beta-lactamase
295 Bla_{TEM-1b}, (Table 2, Table S2 and Table S3). E562 harbours three antibiotic resistance
296 genes, *ampC* located in the chromosome and the *tet(A)* and *bla*_{TEM-1b} genes on an FII
297 plasmid (pAvM_E562_27). The *mer* operon derived from Tn21 is also present in the
298 region (Table 2 and Table S3). The phenotypic resistance profile of E562 matches the
299 genomic profile with resistance to tetracycline (Tet), ampicillin (Amp), amoxicillin-
300 clavulanic acid (Amc) and oxacillin (Oxa) (Table S2). The plasmid pAvM_E36_13
301 contains a complete copy of Tn10, which encodes the *tet(B)*, tetracycline resistance
302 module. Although the AvM_E1373_29 phage-like plasmid is cryptic, related plasmids
303 such as the pHMC2-family of phage-like plasmids [49] (described below), can harbour
304 resistance genes such as *bla*_{CTX-M-14} [50] and *bla*_{CTX-M-15} [51,52].

305

306 Phenotypic intermediate resistance to ampicillin was found in E36 and E1779 encoded
307 by chromosomal *ampC*. Higher MIC values against ampicillin are found in E2980 and
308 E562 strains carrying *bla*_{TEM} genes. Phenotypic resistance to ceftazidime (Caz) and
309 ceftriaxone (Cro) was not found in the isolates, which were consistent with the absence
310 of extended-spectrum beta-lactamase (ESBL) resistance genes in the sequence data.

311

312 Resistance to chloramphenicol (Cm) was found in five isolates, but none of the resistant
313 isolates contained known resistance genes suggesting that chromosomal mutations or
314 presence of efflux pumps may account for this reduced susceptibility.

315 The ETEC reference strains contain several efflux systems which could explain why
316 the genotypic and phenotypic antibiotic resistance profile did not match for all
317 antibiotics. All of the isolates harbour multiple efflux pumps located on the
318 chromosome and plasmids (Table S3 and Figure S12). In E925, a non-synonymous
319 mutation in *acrF* was identified (G1979A) resulting in a substitution from arginine to
320 glutamine (A360Q). The effect on the expression and/or function of the AcrEF efflux
321 pump was not verified.

322

323 Phenotypic resistance to norfloxacin (Nor) was found in 6 of the isolates. The isolates
324 were analysed for chromosomal mutations likely to confer quinolone resistance, using
325 ResFinder but mutations in *gyrA* were only found in one strain, E2980, at position S83A
326 which may confer resistance to nalidixic acid, norfloxacin and ciprofloxacin. However,
327 E2980 was sensitive to Nalidixic acid. Both mutation(s) that alter the target (*gyrA* and
328 *parC*), as well as the presence of efflux pumps, can confer resistance to
329 fluoroquinolones. The majority of the isolates are moderately resistant to Norfloxacin
330 (and Nalidixic acid), both quinolones, which is most likely due to the presence of two
331 efflux pumps, AcrAB-R and AcrEF-R, as only one mutation was identified in *gyrA* of
332 isolate E2980 where usually at least two or more mutations are needed to confer
333 augmented resistance [53].

334

335 ***Identification of phage-like plasmids in ETEC***

336 Two of the ETEC reference strains (E1649 and E1373) harboured phage-like plasmids
337 (pAvM_E1649_9 and pAvM_E1373_29) which encode for DNA metabolism, DNA
338 biosynthesis as well as structural bacteriophage genes (capsid, tail etc.). Both
339 pAvM_E1649_9 and pAvM_E1373_29 contain genes associated with plasmid
340 replication, division and maintenance (i.e. *repA* and *parAB*). Phage-like plasmids are
341 found in various bacterial species, such as *E. coli*, *Klebsiella pneumoniae*, *Yersinia*
342 *pestis*, *Salmonella enterica* serovar Typhi, *Salmonella enterica* serovar Typhimurium,
343 *Salmonella enterica* serovar Derby and *Acinetobacter baumannii* [54]. The plasmid
344 pAvM_E1649_9 belong to the P1 phage-like plasmid family (Figure 2a and Figure
345 S14a) while pAvM_E1373_29 belongs to the pHCM2-family (Figure 2b and Figure
346 S14b) that can be traced back to a likely phage origin similar to the *Salmonella* phage,
347 SSU5 [49]. Both phage-plasmids thus contain replication and/or partition genes of
348 plasmid origin and a complete set of genes that are phage related in function and
349 properties (Figure 2 and Figure S14). Significantly, phage-like plasmid
350 pAvM_E1373_29 falls more within the *E. coli* lineage of pHCM2 phage-like plasmid
351 rather than those found in *Salmonella* species. This indicates that phage-like plasmids
352 have diversified within the bacterial species they were isolated.

353

354

355

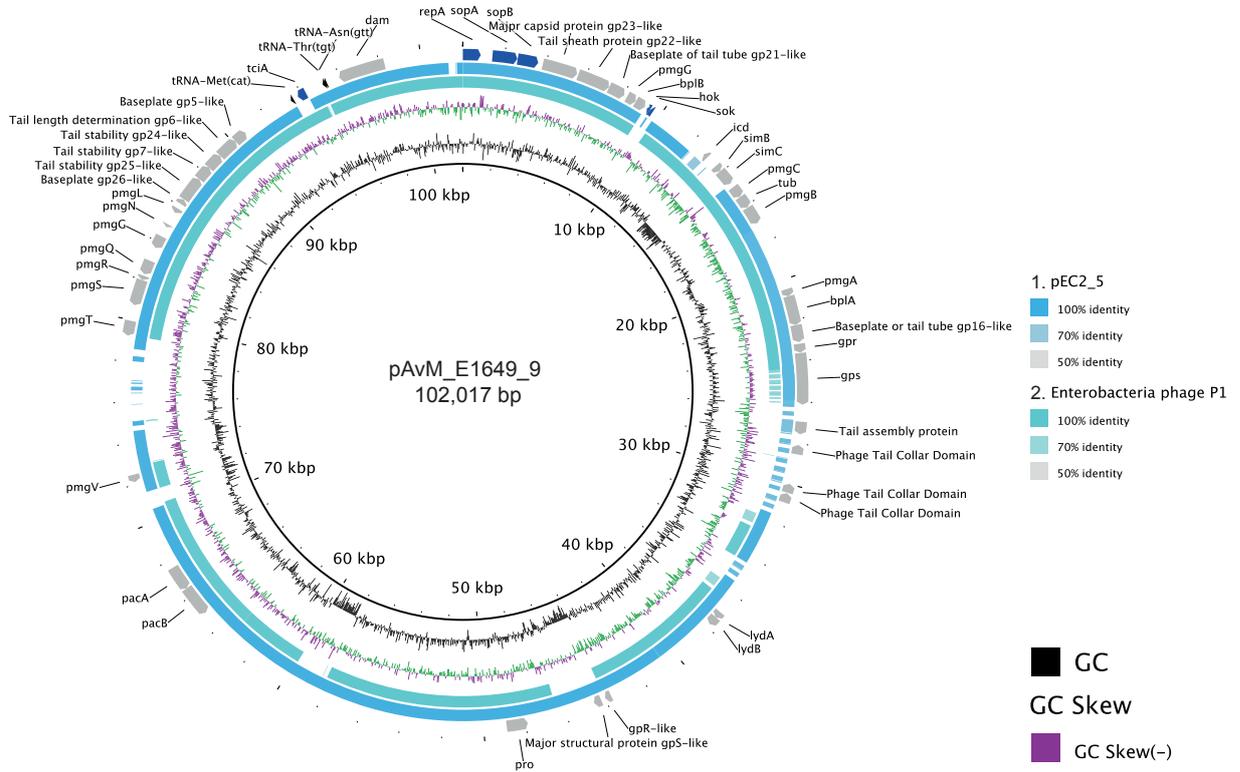
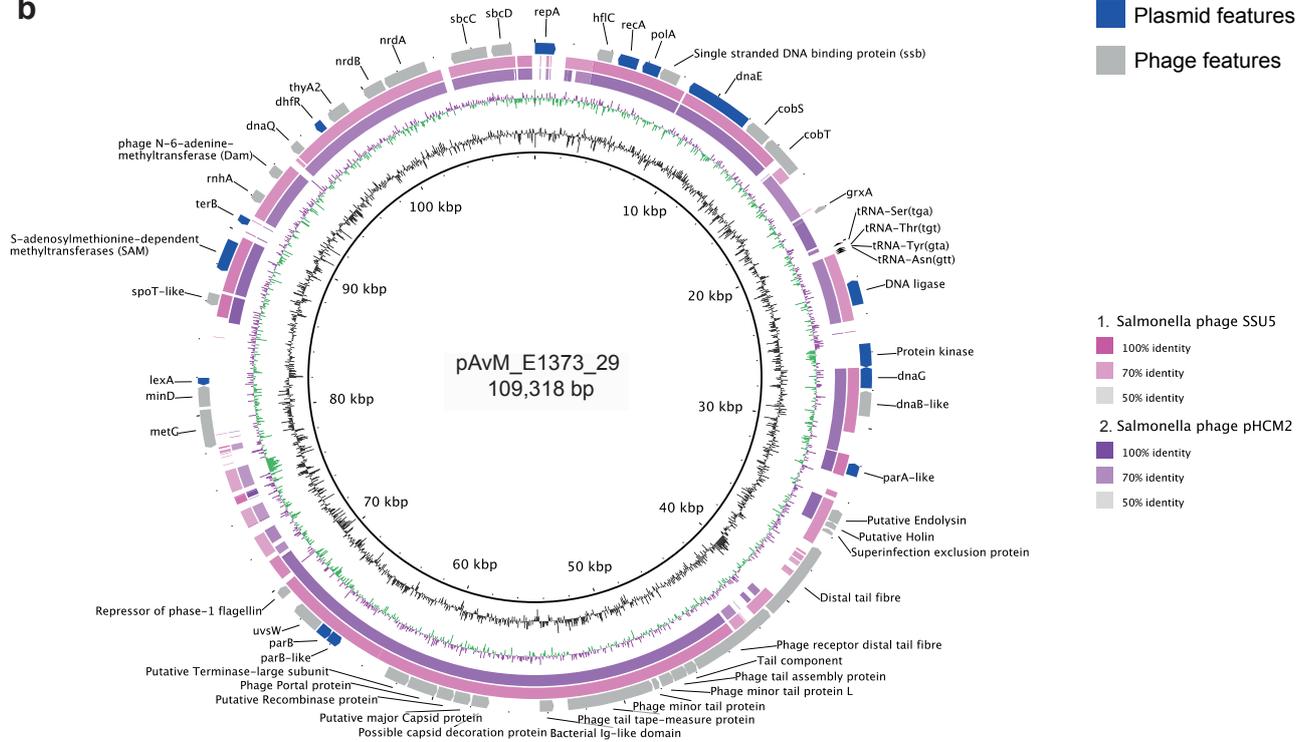
356

357

358

359

360

a**b**

361 **Figure 2:** Comparisons between the identified ETEC phage-plasmids and other similar
362 phage-plasmids using blastn. **a)** pAvM_E1649_9 is a P1-like phage-plasmid here
363 compared to Enterobacteria phage P1 (Escherichia virus P1; NC_005856.1) and
364 pEC2_5 (*E. coli* strain EC2_5; CP041960.1). **b)** pAvM_E2980_29, a phage-plasmid
365 similar to the pHCM2 (Salmonella typhi strain CT18; AL513384.1) and SSU5
366 (Salmonella phage; JQ965645.1) in Salmonella typhi. Blastn comparisons were made
367 using BRIG [84] (v0.95) with the thresholds indicated to the right of each plasmid
368 comparison. Selected phage and plasmid annotations are shown in the outer ring.
369

370

371 Blastn searches confirmed high similarity (at least 80% at the DNA across much of
372 the sequence) of pAvM_E1373_29 to several phage-like plasmids found in *E. coli*
373 including ETEC O169:H41 isolate F8111-1SC3 [20,55], several *bla*_{CTX-M-15} positive
374 phage-like plasmids (pANCO1, pANCO2 [52] and PV234a), as well as a plasmid found
375 in *E. coli* ST648 from wastewater and ST131 isolate SC367ECC [56]. The P1 phage-
376 like plasmid pAvM_1649_9 is most similar to p1107-99K, pEC2_5 isolated from
377 human urine and p2448-3 from a UPEC ST131 isolate isolated from blood. The
378 similarity is most pronounced at the amino acid level. Conservation and synteny are
379 evident when pAvM_1649_9 is compared to P1 phage.

380

381 ***Prophages present and their cargo genes***

382 Prophages may insert into chromosomes and bring along genes required for lysogeny
383 and lytic cycles and cargo genes that are often picked up when DNA is compacted into
384 the capsid. Cargo genes can significantly benefit the host bacterium by providing
385 additional elements to defence against phage or immune evasion and finally,
386 environmental survival. PHASTER analyses identified prophages in the chromosomes
387 of all ETEC reference isolates and some of the plasmids (Table S4). Putative tellurite
388 resistance operons in isolates E925, E36, E2980 and E1373 were all located in

389 prophages. In addition, *eata* (in E925, and E1649) and *estA* (STh) genes (E36) were
390 prophage cargo genes.

391

392 Many prophage cargo genes identified in this study have properties related to inhibition
393 of cell division. Among these are a variety of *kil* genes which can enhance host bacterial
394 survival in the presence of some antibiotics [57]. Some genes that are core entities
395 within many prophages, such as *zapA* (from E1779_Pph_6), *dicB* and *dicC* (found in
396 phage E1779_Pph_7), also have similar effects as they can inhibit cell division in the
397 presence of antibiotics which raise the broader question in terms of how they are
398 beneficial to the host bacterium.

399

400 A different gene of interest is the *yfdR* gene identified in E1779_Pph_7 (gene
401 E1779_04412). YfdR curtails cellular division by inhibiting DNA replication under
402 stress conditions encountered by the bacterial cell. Similarly, the *iram* gene located in
403 phage E1441_Pph_2 plays a role in RpoS stability.

404

405 OmpX homologs were found in numerous phages in this study. They are trans-
406 membrane located and play a role in virulence as well as antibiotic resistance [58]. PerC
407 is often associated with EPEC plasmids, where it seems to have a regulatory role for
408 the attaching and effacing gene, *eaeA* [59]. Its presence as cargo within an ETEC strain
409 phage, E1779_Pph_7 located on the chromosome, is intriguing. Its ability to regulate
410 other virulence genes is yet to be determined. Within the same phage, a *gntR*-like
411 regulatory gene was identified. This gene plays a role in gluconate utilisation and
412 induction of the Entner-Doudoroff pathway [60].

413

414 **Discussion**

415 ETEC strains have previously been shown to fall into globally spread genetically
416 conserved lineages which encompass strains with specific virulence factor profiles [18].
417 The currently widely used ETEC reference strains H10407 (CFA/I) and E24377A
418 (CS1+CS3) are highly divergent from other strains with the same virulence profile
419 sequenced more recently [18] and highlights the need for relevant and representative
420 ETEC reference strains and genomes. The long-read sequenced isolates presented here
421 comprise complete reference genomes with separate chromosomal and plasmid
422 sequences that allow more detailed studies of ETEC and *E. coli* phylogeny. The
423 reference strains are representative isolates of their respective lineage and cluster
424 phylogenetically together with different ETEC isolates sequenced by several other
425 groups (Figure S10).

426 Previous studies confirmed that ETEC belongs to lineages that have spread globally.
427 These analyses were mainly dependent on the shared core genome of chromosomal
428 genes while conservation of plasmids was indicated by the association between the
429 plasmid-borne toxin and CFs and lineage [18]. Analysis of the plasmids sequenced in
430 the present study showed that the conservation within ETEC lineages also include
431 plasmids.

432 Blast analyses confirm that the plasmids identified in this study are often highly
433 homologous to other plasmids identified by either long-read or Illumina sequencing
434 present in GenBank. For instance, the 94.5 kb plasmid pAvM_1441_18 was 98%
435 identical to two 96 kb and 82 kb plasmids belonging to ETEC O25: H16 isolates ATCC
436 43886/E2539C1 and 2014EL-1346-6 sequenced by PacBio by Smith *et al.*, [20] (Figure
437 1e and Figure S6). Plasmid pAvM_E1441_18 is the major virulence plasmid of this
438 lineage carrying genes encoding LT and CS6.

439 The larger plasmid in E1441 (pAvM_E1441_17) carries both the genes for ETEC CF
440 CS21 and antibiotic resistance determinants. Furthermore, complete conjugation
441 machinery was present suggesting that this is most likely a self-transmissible plasmid,
442 though this was not confirmed. Movement of such a plasmid would result in the spread
443 of ETEC virulence genes and AMR determinants.

444 Interestingly, Wachsmuth et al., [44] analysed transfer frequencies in ETEC O25:H16
445 isolates (the same serogroup was identified in E1441) and found evidence that
446 resistance to tetracycline and sulfathiazole was transferred but not the genes encoding
447 LT [44]. The same study found evidence of two large plasmids of similar size [44]
448 corroborating our findings of two plasmids of similar size in E1441, one with *eltAB* and
449 *cssABCD* without the *tra*-operon (pAvM_E1441_18) and the other putatively mobile
450 plasmid (pAvM_E1441_17) carrying the *sull* and *tet(A)* genes as well as the *lng* operon
451 encoding CF CS21. Since ATCC 43886/E2539C1, E1441 and 2014EL-1346-6, have
452 been isolated in the 1970s, 1997, and 2014, respectively, our findings indicate that
453 E1441 represent an ETEC lineage with stable plasmid content and putative ability to
454 transfer antibiotic resistance and the CS21 operon by transfer of one of the plasmids.
455 Furthermore, pAvM_E1441_17 is a multi-replicon plasmid. Multi-replicon plasmids
456 have been described as a way to broaden their host range, i.e. possibility to be
457 transferred between bacteria of different phylogenetic groups [61,62]. Whether this
458 plasmid type is found in other *E. coli* remains to be investigated but the finding that the
459 L4 lineage retains both plasmids in isolates collected over time and worldwide indicate
460 a strong selective force to keep the extra-chromosomal contents of both plasmids.

461 The ETEC O169:H41 isolate F8111-1SC3 plasmid unnamed 2 [20,55] is highly similar
462 to pAvM_E1373_28 (Figure 1h and Figure S9). The F8111-1SC3 isolate is part of a
463 CDC collection of ETEC isolates from cruise ship outbreaks and diarrheal cases in US

464 1996-2003. The antibiotic resistance profiles of these isolates were determined [55] and
465 most isolates of O group 169 were tetracycline resistant consistent with the findings of
466 the *tet* gene in E1373 isolated in Indonesia in 1996. ETEC diarrhoea caused by
467 O169:H41 and STp CS6 isolates is repeatedly reported to cause diarrhoea, particularly
468 in Latin America [45,63–65]. Among the cruise ship isolates is the sequenced and
469 characterised virulence plasmid pEntYN10 encoding STp and CS6, described as
470 unstable and easily lost in vitro [63,66]. The E1373 plasmid; AvM_E1373_28 is highly
471 homologous to pEntYN10 (Figure 1h and Figure S9) and the virulence profile of ETEC
472 O169: H41 is conserved in isolates collected globally. Hence, the instability of the
473 plasmid is incongruent with current data indicating that plasmids are stable within this
474 lineage and serotype.

475 Interestingly, two distinctive extra-chromosomal elements which are highly similar to
476 P1 and SSU5 phage were identified among the 8 ETEC reference strains sequenced
477 (Figure 2, Figure S14 and Table S4). The SSU5-like element carries several genes that
478 allow it to be functional as a plasmid and belongs to the pHCM2-like family of Phage-
479 Plasmids (Figure 2b) [49]. These plasmids are devoid of virulence factors, transposons
480 and antibiotic markers but, they contain a significant number of DNA metabolism and
481 biosynthesis genes and they may contain bacteriophage inhibitory genes that have not
482 yet been identified. Interestingly, several SSU5 phage-like plasmids have been shown
483 to carry the ESBL gene *bla*_{CTX-M15} in extra-intestinal pathogenic *E. coli* isolates [51].
484 ESBL resistance seems to be absent or low in ETEC and the SSU5 phage-like plasmid
485 pAvM_E1373_29 does not contain antibiotic resistance genes. A recent study
486 investigating the distribution of phage-plasmids show that the phage homologs tend to
487 be more conserved and the plasmid homologs more variable [67]. This is also seen in

488 the phage-plasmids identified here, e.g., genes that could be advantageous to the host
489 cell linked to metabolism and biosynthesis.

490 To summarise, we provide fully assembled chromosomes and plasmids with manually
491 curated annotations that will serve as new ETEC reference genomes. The in-depth
492 analysis of gene content, synteny and correct annotations of plasmids will also help to
493 elucidate other plasmids with and without virulence factors in related bacterial species.
494 The ETEC reference genomes compared to other long-read sequenced ETEC genomes
495 confirm that the major ETEC lineages harbour conserved plasmids that have been
496 associated with their respective background genomes for decades. This confirms that
497 the plasmids and chromosomes of ETEC are both crucial for ETEC virulence and
498 success as pathogens.

499

500 **Methods**

501 *Selection of strains*

502 Initially one to two ETEC strains within each of the lineage (L1-L7)-specific CF profile
503 were chosen from the University of Gothenburg large collection of ETEC strains [18]
504 for PacBio sequencing. The strains were selected based on the location and year of
505 isolation to represent strains isolated from patients with diarrhoea from diverse
506 geographical locations and at different time-points. After the genomes had been
507 sequenced, assembled, circularised and annotated a second selection was made for
508 manual curation of the genomes. This selection was made based on the quality of the
509 genome assembly and the circularisation. The whole genomes of the ETEC reference
510 strains were compared with one or two other long-read sequenced ETEC strains
511 belonging to the same lineage by progressiveMAUVE and showed that the strains are
512 colinear (Figure S15). One representative ETEC genome from each lineage was

513 annotated, with emphasis on the plasmids. The physical ETEC reference strains are
514 available upon request.

515

516 ***Phenotypic toxin and CF analyses***

517 ETEC isolates were identified by culture on MacConkey agar followed by an analysis
518 of LT and ST toxin expression using GM1 ELISAs [43]. The expression of the different
519 CFs was confirmed by dot-blot analysis [43]. Isolates had been kept in glycerol stocks
520 at -70°C , and each strain has been passaged as few times as possible.

521

522 ***Antibiotic susceptibility testing***

523 All ETEC isolates were tested against 14 antimicrobial agents and their minimum
524 inhibitory concentration was determined by broth microdilution using EUCAST
525 methodology [47]. The antimicrobial agents were: ampicillin, amoxicillin-clavulanic,
526 oxacillin, ceftazidime, ceftriaxone, doxycycline, tetracycline, nalidixic acid,
527 norfloxacin, azithromycin, erythromycin, chloramphenicol, nitrofurantoin and
528 sulfamethoxazole-trimethoprim. All antibiotics were purchased from Sigma-Aldrich.
529 The *E. coli* ATCC 25922 was used as quality control. The MIC was recorded visually
530 as the lowest concentration of antibiotic that completely inhibits growth.

531

532 ***DNA extraction and sequencing***

533 Strains from each lineage (L1-L7) were SMRT-sequenced on the PacBio RSII. A
534 hybrid *de novo* assembly was performed combining the reads from both the SMRT-
535 sequenced and Illumina sequenced strains.

536 For Single-Molecule Real-Time (SMRT) sequencing (Pacific Bioscience) long intact
537 strands of DNA are required. The genomic DNA extraction was performed as follows.

538 Isolates were cultured in CFA broth overnight at 37°C followed by cell suspension in
539 TE buffer (10 mM Tris and 1 mM EDTA pH 8.0) with 25% sucrose (Sigma) followed
540 by lysis using 10 mg/ml lysozyme (in 0.25 Tris pH 8.0) (Roche). Cell membranes were
541 digested with Proteinase K (Roche) and Sarkosyl NL-30 (Sigma) in the presence of
542 EDTA. RNase A (Roche) was added to remove RNA molecules. A phenol-chloroform
543 extraction was performed using a mixture of Phenol:Chloroform:Isoamyl Alcohol
544 (25:24:1) (Sigma) in phase lock tubes (5prime). To precipitate the DNA 2.5 volumes
545 99% ethanol and 0.1 volume 3 M NaAc pH 5.2 was used followed by re-hydration in
546 10 mM Tris pH 8.0. DNA concentration was measured using NanoDrop
547 spectrophotometer (NanoDrop). On average 10 µg for PacBio sequencing. Library
548 preparation for SMRT sequencing was prepared according to the manufacturers'
549 (Pacific Biosciences) protocol. The DNA was stored in E buffer and sequenced at the
550 Wellcome Sanger Institute. Isolates were sequenced with a single SMRTcell using the
551 P6-C4 chemistry, to a target coverage of 40–60X using the PacBio RSII sequencer.

552 *Assembly*

553 The resulting raw sequencing data from SMRT sequencing were *de novo* assembled
554 using the PacBio SMRT analysis pipeline
555 (<https://github.com/PacificBiosciences/SMRT-Analysis>) (v2.3.0) utilising the
556 Hierarchical Genome Assembly Process (HGAP) [68]. For all samples, the unfinished
557 assembly produced a single, non-circular, chromosome plus some small contigs, some
558 of which were plasmids or unresolved assembly variants. Using Circlator [69] (v1.1.0),
559 small self-contained contigs in the unfinished assembly were identified and removed,
560 with the remaining contigs circularised. Quiver [68] was then used to correct errors in
561 the circularised region by mapping corrected reads back to the circularised assembly.
562 As the strains had also been short read sequenced, and this data is of higher base quality,

563 the short reads from the Illumina sequencing were used in combination with the long
564 reads using Unicycler [70] to generate high-quality assemblies.

565

566 Fully circularised chromosomes and plasmids were achieved for the majority of the
567 strains. Cross-validation of the assemblies was performed where two or three strains of
568 a lineage were sequenced (Figure S15). A single assembly from each lineage was
569 chosen to act as the representative reference genome, with priority given to assemblies
570 with the most complete and circularised chromosome and plasmids. In total, one
571 chromosome and 5 out of the 29 plasmids could not be circularised (independent on the
572 two strains that were sequenced initially) out of the 8 selected representative strains.
573 These are indicated in Table 2 and Table S1. Between two and five plasmids were
574 identified in the eight strains. Shorter contigs that could not be assembled properly
575 contained phage genes and are included in the genomes and annotated as prophages
576 Table S4). Socru was used to validate the assembly of the chromosome, they all have
577 biologically valid orientation and order of rRNA operons with a type GS1.0, which is
578 seen in most *E. coli* in the public domain [71].

579

580 ***Phylogenetic tree***

581 The phylogenetic relationship between the ETEC reference genomes to other ETEC
582 and *E. coli* commensals and pathotypes was investigated. The following collections
583 were included: ETEC-362 [18], ECOR [72] and the Horesh collection [73] along with
584 additional ETEC genomes from several studies [20,24,26,27,45,74,75]. The reads of
585 identified ETEC genomes from other studies were downloaded from GenBank and
586 assembled using Velvet. Long-read sequenced ETEC genomes were included in the
587 tree and were not re-assembled. The phylogroup of the ETEC strains was determined

588 using ClermonTyping [32] (v20.03). The virulence profile of the ETEC strains was
589 determined using ARIBA [76] (v2.14.16) with default settings using the custom ETEC
590 virulence database (https://github.com/avonm/ETEC_vir_db). A total of 1,066
591 genomes was included in the phylogenetic tree. The alignment of core genes (n = 2,895)
592 identified by Roary [77] (v3.12.0) was converted to a SNP-only alignment using snp-
593 sites [78]. A phylogenetic tree was produced with IQ-TREE [79] (v1.6.10) using a GTR
594 gamma model (GTR+F+I) optimised using the built-in model test and visualised using
595 the R package ggtree [80].

596

597 ***Gene prediction, annotation and comparative analysis***

598 The final assembly was annotated using Prokka [81] (v1.14.6). The annotations of all
599 plasmids generated by Prokka were manually checked using the genome viewer
600 Artemis [82] and Geneious 2018.11.1.5 (<http://www.geneious.com>) together with
601 blastp. Annotations of known ETEC virulence genes (colonisation factors, toxins, *eatA*
602 and *etpBAC*) were added after blast+ [83] analysis using the reference genes available
603 in the ETEC virulence database (https://github.com/avonm/ETEC_vir_db) and their
604 annotations updated accordingly. The LT and ST alleles were determined according to
605 Joffre et al., (https://github.com/avonm/ETEC_toxin_variants_db) [15,17]. Where
606 required, PFAM domains were searched using jackhammer to back up any identified
607 protein using blastp (<https://www.ebi.ac.uk/Tools/hmmer/search/jackhammer>). Blastn
608 and tblastx were used for plasmid comparison, using both NCBI website or within
609 BLAST Ring Image Generator (BRIG) [84] (v0.95).

610

611 ***Incompatibility groups***

612 Due to the discrepancy in databases two approaches was used to determine the Inc
613 groups of the 25 plasmids. PlasmidFinder was used with a threshold for minimum %
614 identity at 95% and minimum coverage of 60%. The plasmids were further
615 characterised by pMLST [34], except for IncY which are a group of prophages that
616 replicate in a similar manner as autonomous plasmids (Additional File 3). IncB/O/K/Z
617 plasmids were further typed by blastn comparison to the reference B/O (M93062), K
618 (M93063) and Z (M93064) replicons.

619

620 *oriT prediction*

621 The location of the *oriT* in the plasmids, if present, was predicted using oriTFinder [85]
622 with Blast E-value cut-off set to 0.01.

623

624 *Genomic antibiotic resistance profiling*

625 The identification of antibiotic resistance genes, located on both the chromosome and
626 plasmid(s) as well as the presence of efflux pumps and porins known to confer
627 resistance to antibiotics. The results were obtained by running ARIBA [76] using the
628 CARD database [86] with the default settings (minimum 90% sequence identity and no
629 length cut-off). ARIBA combines a mapping/alignment and targeted local assembly
630 approach to identify AMR genes and variants efficiently and accurately from paired
631 sequencing reads. The heatmaps were visualised using Phandango [87]. The presence
632 of chromosomal mutations in *gyrA* and *parC* was determined with ResFinder (v3.2)
633 from the Center of Genomic Epidemiology [88].

634

635 *Virulence gene prediction*

636 The ETEC assemblies from the ETEC-NCBI collection (Additional file 4) were
637 screened using abricate [89] with default settings against the ETEC virulence database
638 (https://github.com/avonm/ETEC_vir_db) for virulence gene (including *eatA* and
639 *etpBAC*) predication. A subset of the isolates in the ETEC-NCBI dataset have
640 previously been analysed for the presence of EatA where a sample with negative PCR
641 but positive western blots were included as positive [74]. Here, only isolates harbouring
642 the *eatA* and *etpBAC* genes are considered positive.

643

644 ***Prophage prediction***

645 The complete FASTA sequence of each ETEC reference genome was searched for
646 phage genes and prophages using PHASTER (phaster.ca) [90]. The identified intact
647 prophages are listed in Table S4. All prophage contained cargo genes but only
648 recognisable genes are stated, not any hypothetical. Additional questionable and not
649 intact prophages were identified but have not been included here. The prophages have
650 been given a specific identifier name and are also annotated as a *mobile_element* in the
651 submitted chromosome and or plasmid(s) of each strain.

652

653 ***Insertion sequences***

654 Insertion sequences in the plasmids as well as surrounding the CS2 loci located on the
655 chromosome of E1649 were annotated using both Galileo AMR software [91] and the
656 ISFinder database [92]. Complete and partial IS elements were annotated (>95%
657 identity with hits in ISFinder) along with the present genes encoding transposases.
658 Three new insertion sequences were detected in this analysis and were submitted to
659 ISFinder as TnEc2, TnEc3 and TnEc4. Transposons and other mobile elements

660 (integrons and group II introns) were also identified using Galileo AMR and blastn
661 against public databases.

662

663

664

665 **References**

666 1. Khalil IA, Troeger C, Blacker BF, Rao PC, Brown A, Atherly DE, et al. Morbidity
667 and mortality due to shigella and enterotoxigenic *Escherichia coli* diarrhoea: the
668 Global Burden of Disease Study 1990-2016. *The Lancet Infectious diseases*.
669 2018;18:1229–40.

670 2. Baron S, Evans DJ, Evans DG. *Escherichia Coli* in Diarrheal Disease. 1996;

671 3. Svennerholm A, Lundgren A. Recent progress toward an enterotoxigenic
672 *Escherichia coli* vaccine. *Expert Rev Vaccines*. 2014;11:495–507.

673 4. Svennerholm A-M, Lundgren A. Recent progress toward an enterotoxigenic
674 *Escherichia coli* vaccine. *Expert Review of Vaccines*. 2012;11:495–507.

675 5. Lundgren A, Bourgeois L, Carlin N, Clements J, Gustafsson B, Hartford M, et al.
676 Safety and immunogenicity of an improved oral inactivated multivalent
677 enterotoxigenic *Escherichia coli* (ETEC) vaccine administered alone and together
678 with dmLT adjuvant in a double-blind, randomized, placebo-controlled Phase I study.
679 *Vaccine*. 2014;32:7077–84.

680 6. Harro C, Bourgeois AL, Sack D, Walker R, DeNearing B, Brubaker J, et al. Live
681 attenuated enterotoxigenic *Escherichia coli* (ETEC) vaccine with dmLT adjuvant
682 protects human volunteers against virulent experimental ETEC challenge. *Vaccine*.
683 2019;37:1978–86.

684 7. Qadri F, Svennerholm A-M, Faruque AS, Sack RB. Enterotoxigenic *Escherichia*
685 *coli* in developing countries: epidemiology, microbiology, clinical features, treatment,
686 and prevention. *Clinical microbiology reviews*. 2005;18:465–83.

687 8. von Mentzer A, Tobias J, Wiklund G, Nordqvist S, Aslett M, Dougan G, et al.
688 Identification and characterization of the novel colonization factor CS30 based on
689 whole genome sequencing in enterotoxigenic *Escherichia coli* (ETEC). *Scientific*
690 *reports*. 2017;7:465.

691 9. Gaastra W, Svennerholm A-M. Colonization factors of human enterotoxigenic
692 *Escherichia coli* (ETEC). *Trends Microbiol*. 1996;4:444–52.

693 10. Nada RA, Shaheen HI, Khalil SB, Mansour A, El-Sayed N, Touni I, et al.
694 Discovery and phylogenetic analysis of novel members of class b enterotoxigenic

- 695 *Escherichia coli* adhesive fimbriae. *Journal of clinical microbiology*. 2011;49:1403–
696 10.
- 697 11. Cádiz L, Torres A, Valdés R, Vera G, Gutiérrez D, Levine MM, et al. Coli
698 Surface Antigen 26 Acts as an Adherence Determinant of Enterotoxigenic
699 *Escherichia coli* and Is Cross-Recognized by Anti-CS20 Antibodies. *Frontiers in*
700 *microbiology*. 2018;9:248.
- 701 12. Canto FD, O’Ryan M, Pardo M, Torres A, Gutiérrez D, Cádiz L, et al. Chaperone-
702 Usher Pili Loci of Colonization Factor-Negative Human Enterotoxigenic *Escherichia*
703 *coli*. *Frontiers in Cellular and Infection Microbiology*. 2017;6:CD009029.
- 704 13. Grewal HM, Valvatne H, Bhan MK, Dijk L van, Gaastra W, Sommerfelt H. A
705 new putative fimbrial colonization factor, CS19, of human enterotoxigenic
706 *Escherichia coli*. *Infect Immun*. 1997;65:507–13.
- 707 14. Pichel M, Binsztein N, Viboud G. CS22, a novel human enterotoxigenic
708 *Escherichia coli* adhesin, is related to CS15. *Infection and immunity*. 2000;68:3280–
709 5.
- 710 15. Joffre E, von Mentzer A, Ghany MAE, Oezguen N, Savidge T, Dougan G, et al.
711 Allele variants of enterotoxigenic *Escherichia coli* heat-labile toxin are globally
712 transmitted and associated with colonization factors. Christie PJ, editor. *Journal of*
713 *bacteriology*. 2015;197:392–403.
- 714 16. Bolin I, Wiklund G, Qadri F, Torres O, Bourgeois AL, Savarino S, et al.
715 Enterotoxigenic *Escherichia coli* with STh and STp genotypes is associated with
716 diarrhea both in children in areas of endemicity and in travelers. *Journal of clinical*
717 *microbiology*. 2006;44:3872–7.
- 718 17. Joffre E, von Mentzer A, Svennerholm A-M, Sjöling A. Identification of new
719 heat-stable (STa) enterotoxin allele variants produced by human enterotoxigenic
720 *Escherichia coli* (EPEC). *International journal of medical microbiology : IJMM*.
721 2016;306:586–94.
- 722 18. von Mentzer A, Connor TR, Wieler LH, Semmler T, Iguchi A, Thomson NR, et
723 al. Identification of enterotoxigenic *Escherichia coli* (EPEC) clades with long-term
724 global distribution. *Nature genetics*. 2014;46:1321–6.
- 725 19. Crossman LC, Chaudhuri RR, Beatson SA, Wells TJ, Desvaux M, Cunningham
726 AF, et al. A commensal gone bad: complete genome sequence of the prototypical
727 enterotoxigenic *Escherichia coli* strain H10407. *Journal of bacteriology*.
728 2010;192:5822–31.
- 729 20. Smith P, Lindsey RL, Rowe LA, Batra D, Stripling D, Garcia-Toledo L, et al.
730 High-Quality Whole-Genome Sequences for 21 Enterotoxigenic *Escherichia coli*
731 Strains Generated with PacBio Sequencing. *Genome announcements*. 2018;6:6167.

- 732 21. Rasko DA, Rosovitz MJ, Myers GS, Mongodin EF, Fricke WF, Gajer P, et al. The
733 pangenome structure of *Escherichia coli*: comparative genomic analysis of *E. coli*
734 commensal and pathogenic isolates. *Journal of bacteriology*. 2008;190:6881–93.
- 735 22. Besser J, Carleton HA, Gerner-Smidt P, Lindsey RL, Trees E. Next-generation
736 sequencing technologies and their application to the study and control of bacterial
737 infections. *Clinical microbiology and infection : the official publication of the*
738 *European Society of Clinical Microbiology and Infectious Diseases*. 2018;24:335–41.
- 739 23. Quainoo S, Coolen JPM, Hijum SAFT van, Huynen MA, Melchers WJG, Schaik
740 W van, et al. Whole-Genome Sequencing of Bacterial Pathogens: the Future of
741 Nosocomial Outbreak Analysis. *Clinical microbiology reviews*. 2017;30:1015–63.
- 742 24. Sahl JW, Steinsland H, Redman JC, Angiuoli SV, Nataro JP, Sommerfelt H, et al.
743 A comparative genomic analysis of diverse clonal types of enterotoxigenic
744 *Escherichia coli* reveals pathovar-specific conservation. Payne SM, editor. *Infection*
745 *and immunity*. 2011;79:950–60.
- 746 25. Sahl JW, Rasko DA. Analysis of global transcriptional profiles of enterotoxigenic
747 *Escherichia coli* isolate E24377A. *Infection and immunity*. 2012;80:1232–42.
- 748 26. Sahl JW, Sistrunk JR, Fraser CM, Hine E, Baby N, Begum Y, et al. Examination
749 of the Enterotoxigenic *Escherichia coli* Population Structure during Human Infection.
750 *mBio*. 2015;6:e00501.
- 751 27. Sahl JW, Sistrunk JR, Baby NI, Begum Y, Luo Q, Sheikh A, et al. Insights into
752 enterotoxigenic *Escherichia coli* diversity in Bangladesh utilizing genomic
753 epidemiology. *Scientific reports*. 2017;7:3402.
- 754 28. Begum YA, Rydberg HA, Thorell K, Kwak Y-K, Sun L, Joffre E, et al. In
755 SituAnalyses Directly in Diarrheal Stool Reveal Large Variations in Bacterial Load
756 and Active Toxin Expression of Enterotoxigenic *Escherichiacoli*and *Vibrio cholerae*.
757 Limbago BM, editor. *mSphere*. 2018;3:e00517-17.
- 758 29. Li B, Sun J, Han L, Huang X, Fu Q, Ni Y. Phylogenetic groups and pathogenicity
759 island markers in fecal *Escherichia coli* isolates from asymptomatic humans in China.
760 *Applied and environmental microbiology*. 2010;76:6698–700.
- 761 30. Tenaillon O, Skurnik D, Picard B, Denamur E. The population genetics of
762 commensal *Escherichia coli*. *Nature reviews microbiology*. 2010;8:207–17.
- 763 31. Clermont O, Bonacorsi S, Bingen E. Rapid and simple determination of the
764 *Escherichia coli* phylogenetic group. *Applied and environmental microbiology*.
765 2000;66:4555–8.
- 766 32. Beghain J, Bridier-Nahmias A, Nagard HL, Denamur E, Clermont O.
767 ClermonTyping: an easy-to-use and accurate in silico method for *Escherichia* genus
768 strain phylotyping. *Microbial genomics*. 2018;4:690.

- 769 33. McArthur AG, Waglechner N, Nizam F, Yan A, Azad MA, Baylay AJ, et al. The
770 comprehensive antibiotic resistance database. *Antimicrobial agents and*
771 *chemotherapy*. 2013;57:3348–57.
- 772 34. Carattoli A, Zankari E, García-Fernández A, Larsen MV, Lund O, Villa L, et al.
773 *PlasmidFinder* and *pMLST*: in silico detection and typing of plasmids. *Antimicrobial*
774 *agents and chemotherapy*. 2014;58:AAC.02412-14-3903.
- 775 35. Kumar P, Luo Q, Vickers TJ, Sheikh A, Lewis WG, Fleckenstein JM. EatA, an
776 Immunogenic Protective Antigen of Enterotoxigenic *Escherichia coli*, Degrades
777 Intestinal Mucin. Payne SM, editor. *Infection and immunity*. 2014;82:500–8.
- 778 36. Patel SK, Dotson J, Allen KP, Fleckenstein JM. Identification and molecular
779 characterization of EatA, an autotransporter protein of enterotoxigenic *Escherichia*
780 *coli*. *Infection and immunity*. 2004;72:1786–94.
- 781 37. Fleckenstein JM, Roy K, Fischer JF, Burkitt M. Identification of a two-partner
782 secretion locus of enterotoxigenic *Escherichia coli*. *Infection and immunity*.
783 2006;74:2245–58.
- 784 38. Roy K, Hilliard GM, Hamilton DJ, Luo J, Ostmann MM, Fleckenstein JM.
785 Enterotoxigenic *Escherichia coli* EtpA mediates adhesion between flagella and host
786 cells. *Nature*. 2009;457:594–8.
- 787 39. Hibberd ML, McConnell MM, Willshaw GA, Smith HR, Rowe B. Positive
788 regulation of colonization factor antigen I (CFA/I) production by enterotoxigenic
789 *Escherichia coli* producing the colonization factors CS5, CS6, CS7, CS17, PCFO9,
790 PCFO159:H4 and PCFO166. *Journal of general microbiology*. 1991;137:1963–70.
- 791 40. Pilonieta MC, Boder MD, Munson GP. CfaD-dependent expression of a novel
792 extracytoplasmic protein from enterotoxigenic *Escherichia coli*. *Journal of*
793 *bacteriology*. 2007;189:5060–7.
- 794 41. Joffre E, Nicklasson M, Álvarez-Carretero S, Xiao X, Sun L, Nookaew I, et al.
795 The bile salt glycocholate induces global changes in gene and protein expression and
796 activates virulence in enterotoxigenic *Escherichia coli*. *Scientific reports*. 2019;9:108.
- 797 42. Nicklasson M, Sjöling Å, von Mentzer A, Qadri F, Svennerholm A-M. Expression
798 of colonization factor CS5 of enterotoxigenic *Escherichia coli* (ETEC) is enhanced in
799 vivo and by the bile component Na glycocholate hydrate. Hensel M, editor. *PLoS*
800 *One*. 2012;7:e35827.
- 801 43. Sjöling Å, Wiklund G, Savarino SJ, Cohen DI, Svennerholm A-M. Comparative
802 analyses of phenotypic and genotypic methods for detection of enterotoxigenic
803 *Escherichia coli* toxins and colonization factors. *Journal of clinical microbiology*.
804 2007;45:3295–301.
- 805 44. Wachsmuth K, Wells J, Shipley P, Ryder R. Heat-labile enterotoxin production in
806 isolates from a shipboard outbreak of human diarrheal illness. *Infect Immun*.
807 1979;24:793–7.

- 808 45. Hazen TH, Nagaraj S, Sen S, Permala-Booth J, Canto FD, Vidal R, et al. Genome
809 and Functional Characterization of Colonization Factor Antigen I- and CS6-Encoding
810 Heat-Stable Enterotoxin-Only Enterotoxigenic *Escherichia coli* Reveals Lineage and
811 Geographic Variation. Overall CM, editor. *mSystems*. 2019;4:209.
- 812 46. Rozwandowicz M, Brouwer MSM, Fischer J, Wagenaar JA, Gonzalez-Zorn B,
813 Guerra B, et al. Plasmids carrying antimicrobial resistance genes in
814 *Enterobacteriaceae*. *The Journal of antimicrobial chemotherapy*. 2018;73:1121–37.
- 815 47. EUCAST EC for AST. Determination of minimum inhibitory concentrations
816 (MICs) of antibacterial agents by broth dilution. John Wiley & Sons, Ltd (10.1111);
817 2003 Aug p. ix–xv.
- 818 48. Waters SH, Rogowsky P, Grinsted J, Altenbuchner J, Schmitt R. The tetracycline
819 resistance determinants of RP1 and Tn1721: nucleotide sequence analysis. *Nucleic
820 acids research*. 1983;11:6089–105.
- 821 49. Octavia S, Sara J, Lan R. Characterization of a large novel phage-like plasmid in
822 *Salmonella enterica* serovar Typhimurium. *FEMS Microbiology Letters*. 2015;
- 823 50. Liu P, Li P, Jiang X, Bi D, Xie Y, Tai C, et al. Complete genome sequence of
824 *Klebsiella pneumoniae* subsp. *pneumoniae* HS11286, a multidrug-resistant strain
825 isolated from human sputum. *Journal of bacteriology*. 2012;194:1841–2.
- 826 51. Falgenhauer L, Yao Y, Fritzenwanker M, Schmiedel J, Imirzalioglu C,
827 Chakraborty T. Complete Genome Sequence of Phage-Like Plasmid pECOH89,
828 Encoding CTX-M-15. *Genome announcements*. 2014;2:2227.
- 829 52. Colavecchio A, Jeukens J, Freschi L, Rheault J-GE, Kukavica-Ibrulj I, Levesque
830 RC, et al. Complete Genome Sequences of Two Phage-Like Plasmids Carrying the
831 CTX-M-15 Extended-Spectrum β -Lactamase Gene. *Genome announcements*.
832 2017;5:90.
- 833 53. Jacoby GA. Mechanisms of Resistance to Quinolones. *Clin Infect Dis*.
834 2005;41:S120–6.
- 835 54. Gilcrease EB, Casjens SR. The genome sequence of *Escherichia coli* tailed phage
836 D6 and the diversity of *Enterobacteriales* circular plasmid prophages. *Virology*.
837 2018;515:203–14.
- 838 55. Beatty ME, Bopp CA, Wells JG, Greene KD, Puhr ND, Mintz ED. Enterotoxin-
839 producing *Escherichia coli* O169:H41, United States. *Emerging infectious diseases*.
840 2004;10:518–21.
- 841 56. Cho S, Gupta SK, McMillan EA, Sharma P, Ramadan H, Jové T, et al. Genomic
842 Analysis of Multidrug-Resistant *Escherichia coli* from Surface Water in Northeast
843 Georgia, United States: Presence of an ST131 Epidemic Strain Containing blaCTX-
844 M-15 on a Phage-Like Plasmid. *Microbial Drug Resistance*. 2019;mdr.2019.0306.

- 845 57. Wang X, Kim Y, Ma Q, Hong SH, Pokusaeva K, Sturino JM, et al. Cryptic
846 prophages help bacteria cope with adverse environments. *Nat Commun.* 2010;1:147.
- 847 58. Hu WS, Lin J-F, Lin Y-H, Chang H-Y. Outer membrane protein STM3031
848 (Ail/OmpX-like protein) plays a key role in the ceftriaxone resistance of *Salmonella*
849 *enterica* serovar Typhimurium. *Antimicrob Agents Ch.* 2009;53:3248–55.
- 850 59. Gómez-Duarte OG, Kaper JB. A plasmid-encoded regulatory region activates
851 chromosomal *eaeA* expression in enteropathogenic *Escherichia coli*. *Infect Immun.*
852 1995;63:1767–76.
- 853 60. Murray EL, Conway T. Multiple Regulators Control Expression of the Entner-
854 Doudoroff Aldolase (Eda) of *Escherichia coli*. *J Bacteriol.* 2005;187:991–1000.
- 855 61. Villa L, García-Fernández A, Fortini D, Carattoli A. Replicon sequence typing of
856 IncF plasmids carrying virulence and resistance determinants. *The Journal of*
857 *antimicrobial chemotherapy.* 2010;65:2518–29.
- 858 62. Osborn AM, Tatley FM da S, Steyn LM, Pickup RW, Saunders JR. Mosaic
859 plasmids and mosaic replicons: evolutionary lessons from the analysis of genetic
860 diversity in IncFII-related replicons. *Microbiology+*. 2000;146:2267–75.
- 861 63. Nishikawa Y, Helander A, OGASAWARA J, MOYER NP, HANAOKA M,
862 HASE A, et al. Epidemiology and properties of heat-stable enterotoxin-producing
863 *Escherichia coli* serotype O169[ϕ H41]. *Epidemiology and infection.* 1998;121:31–
864 42.
- 865 64. Torres OR, González W, Lemus O, Pradesaba RA, Matute JA, Wiklund G, et al.
866 Toxins and virulence factors of enterotoxigenic *Escherichia coli* associated with
867 strains isolated from indigenous children and international visitors to a rural
868 community in Guatemala. *Epidemiol Infect.* 2014;143:1662–71.
- 869 65. Sack DA, Shimko J, Torres O, Bourgeois AL, Francia DS, Gustafsson B, et al.
870 Randomised, double-blind, safety and efficacy of a killed oral vaccine for
871 enterotoxigenic *E. Coli* diarrhoea of travellers to Guatemala and Mexico. *Vaccine.*
872 2007;25:4392–400.
- 873 66. Ban E, Yoshida Y, Wakushima M, Wajima T, Hamabata T, Ichikawa N, et al.
874 Characterization of unstable pEntYN10 from enterotoxigenic *Escherichia coli*
875 (ETEC) O169:H41. *Virulence.* 2015;6:735–44.
- 876 67. Pfeifer E, Sousa JAM de, Touchon M, Rocha EPC. Bacteria have numerous
877 phage-plasmid families with conserved phage and variable plasmid gene repertoires.
878 *Biorxiv.* 2020;2020.11.09.375378.
- 879 68. Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al.
880 Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing
881 data. *Nature methods.* 2013;10:563–9.

- 882 69. Hunt M, Silva ND, Otto TD, Parkhill J, Keane JA, Harris SR. Circlator:
883 automated circularization of genome assemblies using long sequencing reads.
884 *Genome biology*. 2015;16:294.
- 885 70. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome
886 assemblies from short and long sequencing reads. Phillippy AM, editor. *PLoS*
887 *computational biology*. 2017;13:e1005595.
- 888 71. Page AJ, Ainsworth EV, Langridge GC. socru: typing of genome-level order and
889 orientation around ribosomal operons in bacteria. *Microb Genom*. 2020;
- 890 72. Patel IR, Gangiredla J, Mammel MK, Lampel KA, Elkins CA, Lacher DW. Draft
891 Genome Sequences of the Escherichia coli Reference (ECOR) Collection. *Microbiol*
892 *Resour Announc*. 2018;7:e01133-18.
- 893 73. Horesh G, Blackwell GA, Tonkin-Hill G, Corander J, Heinz E, Thomson NR. A
894 comprehensive and high-quality collection of Escherichia coli genomes and their
895 genes. *Microb Genom*. 2021;
- 896 74. Kuhlmann FM, Martin J, Hazen TH, Vickers TJ, Pashos M, Okhuysen PC, et al.
897 Conservation and global distribution of non-canonical antigens in Enterotoxigenic
898 Escherichia coli. Torres AG, editor. *PLoS neglected tropical diseases*.
899 2019;13:e0007825.
- 900 75. Rasko DA, Canto FD, Luo Q, Fleckenstein JM, Vidal R, Hazen TH. Comparative
901 genomic analysis and molecular examination of the diversity of enterotoxigenic
902 Escherichia coli isolates from Chile. *Plos Neglect Trop D*. 2019;13:e0007828.
- 903 76. Hunt M, Mather AE, Sánchez-Busó L, Page AJ, Parkhill J, Keane JA, et al.
904 ARIBA: rapid antimicrobial resistance genotyping directly from sequencing reads.
905 *bioRxiv*. 2017;118000.
- 906 77. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, et al. Roary:
907 Rapid large-scale prokaryote pan genome analysis. *Bioinformatics*. 2015;31:3691–3.
- 908 78. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, et al. SNP-sites:
909 rapid efficient extraction of SNPs from multi-FASTA alignments. *Microbial*
910 *genomics*. 2016;2:e000056.
- 911 79. Nguyen L-T, Schmidt HA, Haeseler A von, Minh BQ. IQ-TREE: A Fast and
912 Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies.
913 *Mol Biol Evol*. 2015;32:268–74.
- 914 80. Yu G. Using ggtree to Visualize Data on Tree-Like Structures. *Curr Protoc*
915 *Bioinform*. 2020;69:e96.
- 916 81. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*.
917 2014;30:2068–9.

- 918 82. Carver T, Harris SR, Berriman M, Parkhill J, McQuillan JA. Artemis: an
919 integrated platform for visualization and analysis of high-throughput sequence-based
920 experimental data. *Bioinform Oxf Engl*. 2011;28:464–9.
- 921 83. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al.
922 BLAST+: architecture and applications. *Bmc Bioinformatics*. 2009;10:421.
- 923 84. Alikhan N-F, Petty NK, Zakour NLB, Beatson SA. BLAST Ring Image Generator
924 (BRIG): simple prokaryote genome comparisons. *Bmc Genomics*. 2011;12:402.
- 925 85. Li X, Xie Y, Liu M, Tai C, Sun J, acids ZDN, et al. oriTfinder: a web-based tool
926 for the identification of origin of transfers in DNA sequences of bacterial mobile
927 genetic elements. *academic.oup.com*.
- 928 86. Jia B, Raphenya AR, Alcock B, Wagglehner N, Guo P, Tsang KK, et al. CARD
929 2017: expansion and model-centric curation of the comprehensive antibiotic
930 resistance database. *Nucleic acids research*. 2017;45:D566–73.
- 931 87. Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, Harris SR.
932 Phandango: an interactive viewer for bacterial population genomics. *Bioinform Oxf*
933 *Engl*. 2017;34:292–3.
- 934 88. Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, et al.
935 Identification of acquired antimicrobial resistance genes. *J Antimicrob Chemother*.
936 2012;67:2640–4.
- 937 89. Seemann T. Abricate [Internet]. Available from:
938 <https://github.com/tseemann/abricate>
- 939 90. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better,
940 faster version of the PHAST phage search tool. *Nucleic acids research*. 2016;44:W16-
941 21.
- 942 91. Partridge SR, Tsafnat G. Automated annotation of mobile antibiotic resistance in
943 Gram-negative bacteria: the Multiple Antibiotic Resistance Annotator (MARA) and
944 database. *J Antimicrob Chemother*. 2018;73:883–90.
- 945 92. Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M. ISfinder: the reference
946 centre for bacterial insertion sequences. *Nucleic Acids Res*. 2006;34:D32–6.
- 947
- 948
- 949
- 950
- 951

952

953

954

955

956

957 **Acknowledgements**

958 No acknowledgements to mention.

959

960 **Authors' contributions**

961 AvM conceived and designed the experiments, performed the experiments, analysed
962 the data, contributed reagents/materials/analysis tools, prepared figures and/or tables,
963 authored the paper and approved the final draft.

964 GB and AvM annotated all IS elements, transposons as well as other mobile elements,
965 contributed to the paper and approved the final draft.

966 DP and AvM annotated the identified prophages, contributed to the paper and approved
967 the final draft.

968 CB performed the *in silico* analysis of the genomic antibiotic resistance profiling,
969 contributed to the paper and approved the final draft.

970 EJ performed the antibiotic resistance profiling, contributed to the paper and approved
971 the final draft.

972 AJP assembled the genomes, contributed to the paper and approved the final draft.

973 AMS conceived and designed the experiments contributed to the paper and approved
974 the final draft.

975 GD conceived and designed the experiments and approved the final draft.

976 ÅS conceived and designed the experiments, analysed data, authored the paper and
977 approved the final draft.

978

979 **Competing interests**

980 The authors declare that there are no conflicts of interest.

981

982 **Data availability**

983 The datasets supporting the conclusions of this article are included within the articles
984 and its additional files. The sequencing data generated in this study has been
985 submitted to EMBL (Additional file 4 and 5). The physical ETEC reference strains
986 can be requested by contacting the corresponding author Astrid von Mentzer
987 (avm@sanger.ac.uk or mentzerv@chalmers.se). The database used for annotating
988 ETEC virulence factors, ETEC virulence database, including the LT and ST alleles
989 can be found in the github repositories: https://github.com/avonm/ETEC_vir_db and
990 https://github.com/avonm/ETEC_toxin_variants_db.

991 An interactive version of the core genome phylogeny of the 1,065 *E. coli* and ETEC
992 isolates along with the ETEC reference strains (Figure S10) reported here is
993 accessible at <https://microreact.org/project/2ZZzaHzeXbMEw9U2MAk7pK?tt=cr>

994

995 **Availability of material**

996 Requests for obtaining clinical isolates collected as part of this study should be
997 addressed to the corresponding author. Exchange of clinical isolates should always be
998 in agreement with the University of Gothenburg.

999

1000 **Ethics declarations**

1001 Not applicable.

1002

1003 **Consent for publication**

1004 Not applicable.

1005

1006 **Additional information**

1007 **Additional file 1:** Supplemental Figures and Tables
1008 **Additional file 2:** Detailed description of ETEC reference plasmids
1009 **Additional file 3:** Excel file with plasmid classification – Inc groups
1010 **Additional file 4:** Excel file with metadata of ETEC and *E. coli* genomes included in
1011 the phylogenetic tree
1012 **Additional file 5:** Excel file with a compiled list of all accession numbers for
1013 chromosomes and plasmids
1014

1015 **Funding**

1016 AvM, AMS and ÅS were supported by the Swedish Foundation for Strategic Research
1017 (grant nr. SB12-0072). AvM was also supported by The Swedish Research Council
1018 (grant nr. 2018-06828) and the Swedish Society for Medical Research (P18-0140). AJP
1019 was supported by the Biotechnology and Biological Sciences Research Council
1020 (BBSRC); this research was funded by the BBSRC Institute Strategic Programme
1021 Microbes in the Food Chain BB/R012504/1. GD was supported by the Wellcome Trust
1022 (grant WT 098051).

1023

1024 **Author details**

1025 ¹Department of Microbiology and Immunology, Sahlgrenska Academy, University of
1026 Gothenburg, Medicinaregatan 7, 413 90 Gothenburg, Sweden

1027 ²Wellcome Sanger Institute, Hinxton, CB10 1SA, Cambridge, UK

1028 ³EMBL-EBI, Hinxton, CB10 1SA, Cambridge, UK

1029 ⁴University of Cambridge Department of Medicine, Addenbrookes Hospital

1030 Box 157, Hills Rd, CB2 0QQ, Cambridge, UK

- 1031 ⁵ Department of Microbiology, Tumor and Cell Biology, Biomedicum, Karolinska
1032 Institutet, Solnavägen 9, 171 65 Solna, Stockholm, Sweden
- 1033 ⁶Quadram Institute Bioscience, Norwich Research Park, NR4 7UQ Norwich, UK
- 1034 [#]Currently affiliated with Chalmers University of Technology, Chalmers Tvärgata 3,
1035 413 58 Gothenburg, Sweden
- 1036

Figures

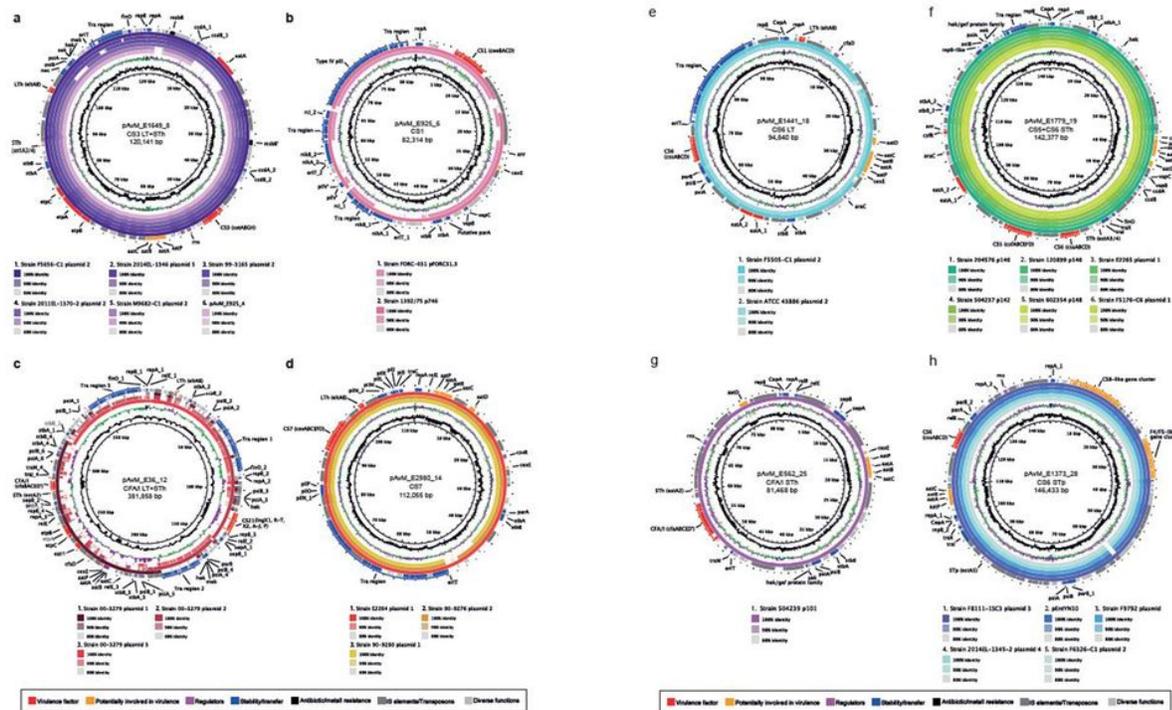


Figure 1

Comparison between the ETEC reference plasmids harbouring colonisation factors and other PacBio-sequenced ETEC plasmids using blastn. a) pAvM_E1649_8 (CS3) as reference and pAvM_925_4 (CS3) compared to the following ETEC plasmids: F5656-C1 plasmid 2 (USA, CP024262.1), 2014-EL-1346-6 plasmid 5 (2014, USA; CP024237.1), 99-3165 plasmid 2 (USA; CP029980.1), 2011EL-1370-2 plasmid 2 (2011, USA; CP022914.1) and M9682-C1 plasmid 2 (1975, USA; CP024277.1). b) pAvM_E925_6 (CS1) compared to ETEC plasmids pFORC31.3 (2004, Korea; CP013193.1) and 1392/75 p746 (1973; FN822748.1). c) pAvM_E36_12 (CFA/I) compared to plasmids 1-3 (p1: CP024294.1; p2 CP024295.1; p3: CP024296.1) from ETEC strain 00-3279 (USA). d) pAvM_E2980_14 (CS7) compared to E2264 plasmid 1 (2006, Bangladesh; CP023350.1), 90-9276 plasmid 2 (1988, Bangladesh; CP024298.1) and 90-9280 plasmid 1 (1988, Bangladesh; CP024241.1). e) pAvM_E1441_18 (CS6) compared to F5505-C1 plasmid 2 (2013, Sweden; CP023259.1) and ATCC 43886 plasmid 2 (CP024255.1). f) pAvM_E1779_19 (CS5+CS6) compared to 204576 p146 (2010, Mali; CP025908.1), 120899 p146 (2012, Gambia; CP025917.1), E2265 plasmid 1 (2006, Bangladesh; CP023347.1), 504237 p142 (2010, India; CP025863.1), 602354 p148 (2009, Bangladesh; CP025848.1) and F5176-C6 plasmid 1 (1997; CP024668.1). g) pAvM_E562_25 (CFA/I) compared to p504239_101 (2010, India; CP025860.1). h) pAvM_E1373_28 (CS6) compared to F8111-1SC3 plasmid 3 (USA; CP024272.1), pEntYN10 (1991, Japan; AP014654.2), F9792 plasmid (USA; CP023274.1), 2014EL-1345-2 plasmid 4 (2014, USA; CP024227.1) and F6326-C1 plasmid 2 (1998, USA; CP024265.1). The thresholds chosen for the blastn are shown in the key below each plasmid comparison. The colour code for the annotations are listed at the bottom of the figure. The two most

inner rings depict GC content in black and GC Skew- in purple and GC Skew+ in green. The figures were generated using BRIG [84] v0.95.

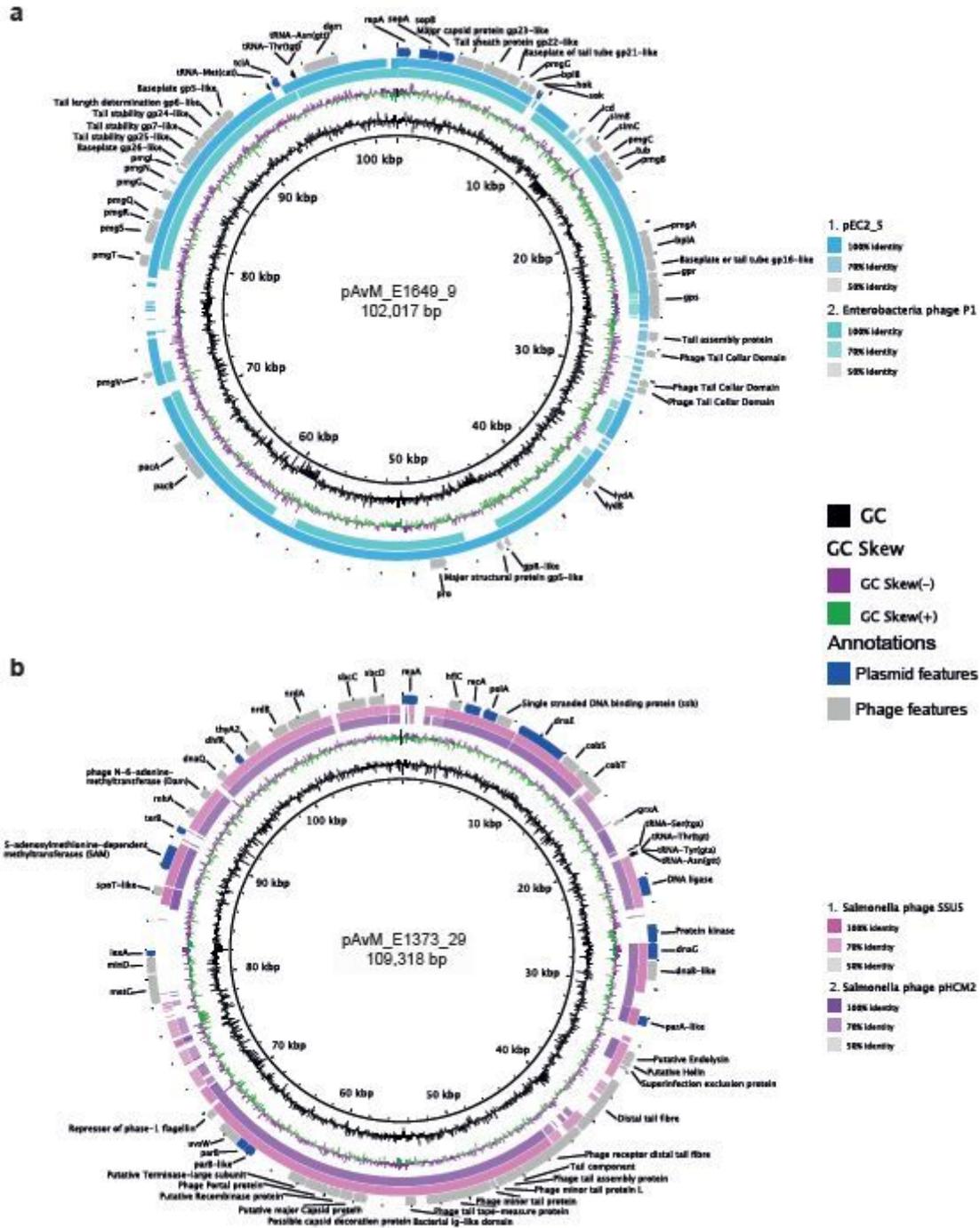


Figure 2

Comparisons between the identified ETEC phage-plasmids 361 and other similar phage-plasmids using blastn. a) pAvM_E1649_9 is a P1-like phage-plasmid here compared to Enterobacteria phage P1 (Escherichia virus P1; NC_005856.1) and pEC2_5 (E. coli strain EC2_5; CP041960.1). b) pAvM_E2980_29,

a phage-plasmid similar to the pHCM2 (*Salmonella typhi* strain CT18; AL513384.1) and SSU5 (*Salmonella* phage; JQ965645.1) in *Salmonella typhi*. Blastn comparisons were made using BRIG [84] (v0.95) with the thresholds indicated to the right of each plasmid comparison. Selected phage and plasmid annotations are shown in the outer ring.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Table1.pdf](#)
- [Table2.pdf](#)
- [Additionalfile1.pdf](#)
- [AdditionalFile2finalPlasmiddescription.pdf](#)
- [Additionalfile3Incprofiles.xlsx](#)
- [Additionalfile4metadatav2final.xlsx](#)
- [Additionalfile5accno.xlsx](#)