**A commensal gone bad: complete genome sequence of the prototypical enterotoxigenic *Escherichia coli* strain H10407**

Lisa C. Crossman[1†§], Roy R. Chaudhuri[2§], Scott A. Beatson[3], Timothy J. Wells[4], Mickael Desvaux[4ƒ], Adam F. Cunningham[4], Nicola K. Petty[1], Vivienne Mahon[5], Carl Brinkley[6], Jon L. Hobman[7], Stephen J. Savarino[6], Susan M. Turner[4], Mark J. Pallen[8], Charles W. Penn[8], Julian Parkhill[1], A. Keith Turner[1], Timothy J. Johnson[9], Nicholas R. Thomson[1], Stephen G.J. Smith[5], Ian R. Henderson[4*]

[1]The Wellcome Trust Sanger Institute, Genome Campus, Hinxton, Cambridge, United Kingdom; [2]Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom; [3]School of Chemistry and Molecular Biosciences, University of Queensland, Brisbane, Australia; [4]School of Immunity and Infection and [8]School of Biosciences, University of Birmingham, Birmingham, United Kingdom; [5]Department of Clinical Microbiology, School of Medicine, Trinity College Dublin, Dublin, Ireland; [6]Department of Enteric Infections, Division of Communicable Diseases and Immunology, Walter Reed Army Institute of Research, Silver Spring, Maryland, USA; [7]School of Biosciences, The University of Nottingham, Sutton Bonington, United Kingdom; [9]Department of Veterinary and Biomedical Sciences, University of Minnesota, Saint Paul, Minnesota.

*Corresponding author. Mailing address: School of Immunity and Infection, University of Birmingham, Birmingham, B15 2TT, United Kingdom. Phone: +44 121 4144368. Email: i.r.henderson@bham.ac.uk

Present address: [†] The Genome Analysis Centre, Norwich, United Kingdom and [ƒ]INRA, UR454 Microbiology, F-63122 Saint-Genès Champanelle, France

[§]LCC and RRC contributed equally to this investigation.

Running Title: **ETEC genome sequence**

**ABSTRACT**

**In most cases *Escherichia coli* exists as a harmless commensal organism but may on occasion cause intestinal and/or extraintestinal disease. Enterotoxigenic *E. coli* are the predominant cause of *E. coli*-mediated diarrhea in the developing world and are responsible for a significant portion of paediatric deaths. In this study we determined the complete genomic sequence of *E. coli* H10407, a prototypical strain of enterotoxigenic *E. coli*, which reproducibly elicits diarrhea in human volunteer studies. We performed genomic and phylogenetic comparisons with other *E. coli* strains revealing that the chromosome is closely related to the non-pathogenic commensal strain *E. coli* HS and to the laboratory strains *E. coli* K-12 and C. Furthermore, these analyses demonstrated that there were no chromosomally-encoded factors unique to any sequenced ETEC strains. Comparison of the *E. coli* H10407 plasmids with those from several ETEC strains revealed the plasmids had a mosaic structure but that several loci were conserved amongst ETEC strains. This study provides a genetic context for the vast amount of experimental and epidemiological data published thus far.**

**INTRODUCTION**

Current dogma suggests the Gram-negative motile bacterium *Escherichia coli* colonises the infant gut within hours of birth and establishes itself as the predominant facultative anaerobe of the colon for the remainder of life (3, 59). While the majority of *E. coli* strains maintain this harmless existence some strains have adopted a pathogenic lifestyle. Contemporary tenets suggest that pathogenic strains of *E. coli* have acquired genetic elements, which encode virulence factors and enable the

2

55 organism to cause disease (12). The large repertoire of virulence factors enables *E.*

56 *coli* to cause a variety of clinical manifestations including intestinal infections

57 mediating diarrhea and extraintestinal infections, such as urinary tract infections,

58 septicaemia and meningitis. Based on clinical manifestation of disease, the

59 repertoire of virulence factors, epidemiology and phylogenetic profiles, strains

60 causing intestinal infections can be divided into six separate pathotypes viz.

61 enteroaggregative *E. coli* (EAEC), enteroinvasive (EIEC), enteropathogenic *E. coli*

62 (EPEC), enterohaemorrhagic *E. coli* (EHEC), diffuse adhering *E. coli* (DAEC) and

63 enterotoxigenic *E. coli* (ETEC) (33, 35, 39).

64 ETEC is responsible for the majority of *E. coli*-mediated cases of human diarrhea

65 worldwide. It is particularly prevalent amongst children in developing countries

66 where sanitation and clean supplies of drinking water are inadequate, and in

67 travellers to such regions. It is estimated that there are 200 million incidences of

68 ETEC infection annually resulting in hundreds of thousands of deaths in children

69 under the age of 5 (55, 64). The essential determinants of ETEC virulence are

70 traditionally considered to be colonization of the host small intestinal epithelium via

71 plasmid-encoded colonization factors (CFs), and subsequent release of plasmid-

72 encoded heat-stable (ST) and/or heat-labile (LT) enterotoxins that induce a net

73 secretory state leading to profuse watery diarrhea (20, 62). More recently, additional

74 plasmid-encoded factors have been implicated in the pathogenesis of ETEC, namely

75 the EatA serine protease autotransporter (SPATE) and the EtpA protein that acts as

76 an intermediate in the adhesion between bacterial flagella and host cells (23, 32, 42,

77 46). Furthermore, a number of chromosomal factors are thought to be involved in

78 virulence e.g. the invasin Tia, the TibA adhesin/invasin and LeoA, a GTPase of

79 unknown function (14, 21, 22). *E. coli* H10407 is considered a prototypical ETEC

3

80  strain; it expresses colonization factor antigen 1 (CFA/I) and the heat-stable and heat

81  labile toxins. Loss of a 94.8-kb plasmid encoding CFA/I and a gene for ST

82  enterotoxin from *E. coli* strain H10407 leads to reduced ability to cause diarrhea (17).

83  Here we report the complete genome sequence and virulence factor repertoire of the

84  prototypical ETEC strain H10407, the nucleotide sequence and gene repertoire of

85  the plasmids from ETEC strain E1392/75 and we describe a novel conserved

86  secretion system associated with the sequenced ETEC strains.

87  **MATERIALS AND METHODS**

88  **Bacterial strains and sequencing.**

89  The ETEC O78:H11:K80 strain H10407 was isolated from an adult with cholera-like

90  symptoms in the course of an epidemiologic study in Dacca, Bangladesh prior to

91  1973 (19) and was shown to cause diarrhea in adult volunteers (6, 17). The *E. coli*

92  H10407 isolate that was sequenced was from the Walter Reed Army Institute of

93  Research (WRAIR) cGMP stock manufactured in February 1998 as Lot 0519.  The

94  whole genome was sequenced to a depth of 8 x coverage from pUC19 (insert size

95  2.8-5 kb) and pMAQ1b (insert size 5.5-10 kb) small insert libraries. Sanger

96  Sequencing was carried out using Amersham Big-Dye (Amersham, UK) terminator

97  chemistry on ABI3700 sequencing machines. End sequences from larger insert

98  plasmid (pBACe3.6, 20-30 kb insert size) libraries were used as a scaffold.

99  Sequence reads were assembled into contigs with Phrap (Green P, unpublished)

100  and finished using GAP4 as described previously (33). The plasmids from ETEC

101  O6:H16:K15 strain E1392/75, which was isolated from a patient in Hong Kong with

102  diarrhea, expresses the CFA/II (CS1 and CS3) colonization factors and produces the

103  ST and LT toxins, were also sequenced using a similar approach (7, 50, 60).

104  Plasmid DNA for ETEC E1392/75 was provided by Acambis UK.

4

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

Annotation was carried out using the genome viewer Artemis (47). Coding sequences were predicted using the gene prediction programs Orpheus (26), Glimmer2 (11) and Glimmer3 (10), then manually curated. Protein domains were marked up using Pfam (48) and transmembrane domains and signal sequences were predicted using TMHMM and SignalP, respectively (15, 37). Annotation was transferred from previously annotated *E.coli* genomes to orthologous genes and manually curated. A homologue was considered to be present if a hit was found with >60% identity over at least 80% of the length of the query protein. Regions of difference and plasmids were annotated and curated manually. The annotated genome sequence of ETEC H10407 and the plasmids from ETEC H10407 and E1392/75 have been deposited in the EMBL databases (accession number: FN649414 for ETEC H10407 complete chromosome; see Tables 1 and 2 for general features of the nucleotides sequences and accession numbers for the plasmids).

**RESULTS AND DISCUSSION**

**Structure and general features of ETEC H10407 chromosome.**

The ETEC H10407 genome consists of a circular chromosome of 5,153,435 bp and four plasmids designated pETEC948, pETEC666, pETEC58 and pETEC52, respectively. The general features of the ETEC H10407 chromosome are presented in Table 1 and the plasmids in Table 2. We identified 4746 protein-coding genes (CDSs) in the chromosome, 33 (0.67%) of which do not have any match in the database, 579 (11.67%) encode conserved hypothetical proteins, with no known function and 503 (10.14%) are genes associated with mobile elements such as integrases, transposases, or phage related. We have identified 25 regions of difference (ROD) that occur in the ETEC H10407 genome and are differentially

5

130    distributed among the other sequenced *E. coli* chromosomes (Figure 1; Table S1).

131    The combined size of these RODs is 755,359 bp (14.7% of the chromosome) and

132    includes nine prophages, designated ETP29, 33, 86, 128, 216, 284, 295, 468 and

133    507, where the numeric designations denote their approximate positions (x 10,000

134    bp) on the chromosome. None appeared to carry cargo genes related to virulence.

135    **Comparative genomics of the ETEC H10407 chromosome.**

136    Previously, a phylogeny was constructed based on the concatenated sequences of

137    2,173 genes that are conserved in all *E. coli* strains and in *Escherichia albertii* and

138    *Escherichia fergusonii*, which were included as outgroup sequences (4). The

139    established *E. coli* sub-groups (A, B1, B2, D and E) are all monophyletic with the

140    exception of group D, which is divided at the root. In agreement with previous optical

141    mapping experiments (5), *E. coli* H10407 is located in the A subgroup with the non-

142    pathogenic laboratory strains *E. coli* K-12 and C and the non-pathogenic commensal

143    isolate *E. coli* HS.  The majority of commensal strains of bacteria belong to the A

144    subgroup (59).

145    Comparison of *E. coli* H10407 with the closely related non-pathogenic *E. coli* K-12, C

146    and HS strains reveals these chromosomes are largely colinear (Fig S1) and that *E.*

147    *coli* H10407 chromosome contains 599 CDSs not present in the non-pathogenic

148    strains (Fig. 2 and Table S2).  The majority (528) of these are clustered in the 25

149    RODs and are predicted to encode prophage genes and other mobility factors.

150    Several genes encode previously described loci specifically associated with ETEC

151    virulence viz. *leoA* (ROD 20), *tia* (ROD 20) and *tib* (ROD 13) (13, 14, 22).  Other

152    genes encode loci previously noted in ETEC H10407 including the degenerate ETT2

153    locus (ROD 18) (45), Antigen 43 (ROD 23) (63), a Type 2 protein secretion locus

154    found in many strains of *E. coli* (ROD 19) (4) and the *ecpP* fimbrial gene cluster also

6

155  found in many *E. coli* strains (ROD 1) (4).  Other RODs encode the Sil/Pco efflux

156  system conferring silver/copper resistance (ROD 2), yersiniabactin (ROD 11), and

157  the O78 serotype O-antigen biosynthetic locus (ROD 14).  The *sil* operon is closely

158  related to *sil* from IncH2 plasmid pMG101 (30, 38, 53) and is adjacent to a partially

159  interrupted copper resistance operon similar to *pco* from plasmid pRJ1004 (2). The

160  *sil/pco* locus is flanked by IS element and phage related sequences suggesting

161  horizontal transfer of these genes. The yersiniabactin iron acquisition locus is widely

162  distributed in *E. coli* and other members of the *Enterobacteriaceae* (49).   The

163  remaining *E. coli* H10407-specific CDSs, which are not present on a ROD, and do

164  not encode prophage or mobility factors, encode the H11 flagellin subunits (CDS

165  2029-2033), an additional copy of Antigen 43 (CDS 2119), and several pseudogenes

166  (CDS 427, 1476, 1573). This data largely agrees with previously published

167  subtractive hybridisation studies (5).

168  If a particular protein plays an important role in ETEC-mediated disease then one

169  would expect the gene encoding it to have a wide distribution amongst ETEC strains.

170  To determine if there were any chromosomally-encoded genes specific for ETEC

171  strains, comparisons were made with *E. coli* strains E24377A and B7A, the only

172  other ETEC strains for which genome sequence data is available (44). Unlike *E. coli*

173  H10407 both the *E. coli* strains E24377A and B7A belong to the B1 subgroup of the

174  *E. coli* phylogeny, a subgroup from which many commensals are derived but also a

175  number of pathogens (4, 59).  Comparison of *E. coli* H10407 with the sequenced

176  ETEC strain E24377A revealed the chromosomes are largely colinear (Fig S2). The

177  genome of ETEC B7A is not finished but experience with other *E. coli* genomes and

178  comparison of the 198 finished ETEC B7A contigs suggests that the chromosome is

179  also largely colinear with the other sequenced ETEC genomes (Fig. S2).  Analyses

180    of the gene content of all three strains revealed 3741 genes conserved in all three

181    strains, of which only 188 are not present in the commensal *E. coli* HS (Fig. 2B and

182    Table S3). The 188 genes identified through this comparison included loci encoding

183    xanthine dehydrogenase (CDS 0339-0343), the Mat fimbriae (CDS 0348-0352),

184    conserved proteins of unknown function (CDS 0673-0678), a flavoprotein electron

185    transfer system (CDS 1730-1734), the colanic exopolysaccharide biosynthetic

186    machinery (CDS 2171-2202), the Fec iron citrate uptake system (CDS 3161-3166), a

187    cellulose synthase system (CDS 3776-3779) and a putative sugar utilisation system

188    (CDS 4145-4154) all of which were present in the non-pathogen *E. coli* K-12 and are

189    widely distributed amongst other *E. coli* (data not shown).  The remainder of the 188

190    genes encode prophage or other mobility factors which are predicted to have no role

191    in virulence.   Of the 599 *E. coli* H10407-restricted genes identified through

192    comparisons with the non-pathogenic *E. coli* strains above (Fig. 2A), 47 were

193    conserved amongst the three pathogenic ETEC isolates. However, these genes

194    were all related to mobile elements and no putative virulence factors were identified.

195    Notably, no significant homologues of *leoA*, *tibC*, *tibA* or *tia* were detected in either

196    *E. coli* E24377A or B7A strongly suggesting these genes are not essential for ETEC-

197    mediated disease.  In conclusion, these data agree with previous observations that

198    the chromosome of *E. coli* H10407 is most closely related to non-pathogenic *E. coli*

199    and the factors mediating diarrhea are not chromosomally encoded thereby

200    indicating the essential virulence factors are encoded on the plasmids (61).

201    **Potential virulence genes encoded on the ETEC plasmids**

202    Since chromosomal comparisons revealed that no chromosomal CDS was unique to

203    all three ETEC strains we next examined the CDSs present on the four plasmids of

204    ETEC H10407. The general characteristics of the plasmids are shown in Table 2.

205    The two larger plasmids (pETEC948 and pETEC666) are reminiscent of conjugative

206    plasmids that are often associated with the carriage of virulence factors whereas the

207    two smaller plasmids (pETEC58 and pETEC52) are homologous to mobilisable

208    plasmids frequently encountered in a variety of bacterial species (24, 34). The latter

209    plasmids have been shown to be mobilizable in the presence of IncF and other

210    plasmid transfer systems (51).  The majority of the CDSs on all four plasmids encode

211    plasmid maintenance and transfer functions, pseudogenes, genes of unknown

212    function not predicted to be involved in virulence and transmissible elements (Table

213    2).  An exhaustive list of the genetic content is unwarranted here, as the complete

214    annotation of the plasmids is provided via the EMBL databases. Nevertheless, there

215    are several noteworthy CDSs, described below, which can be termed "cargo" genes

216    and have a known or putative role in pathogenesis.  Thus, analyses revealed *E. coli*

217    H10407 pETEC948 possesses cargo genes encoding the previously described EatA

218    SPATE (*eatA*), heat-stable enterotoxin STa2 (*sta2*), CFA/I fimbriae and associated

219    regulator (*cfaABCD*), Etp two partner secretion system and associated

220    glycosyltransferase (*etpABC*) (Fig. 3) (18, 23, 42, 66). Analyses of the *E. coli*

221    H10407 pETEC666 plasmid revealed it contains the cargo genes encoding the

222    previously described heat stable enterotoxin STa1 (*sta1*) and the two subunits of LT

223    enterotoxin (*eltA* and *eltB*) (Fig 3) (8, 65).  In addition, the plasmids contain several

224    loci not previously associated with ETEC strains. ETEC H10407 pETEC948

225    possesses genes encoding a Type I secretion locus similar to the dispersin secretion

226    locus (*aatA-P*) described for *E. coli* 042 (Fig. 4)(52).  Associated with this locus is a

227    gene encoding CexE, a previously described secreted protein of ETEC (43), which

228    bears homology to the *E. coli* 042 dispersin protein (Fig. 4). Furthermore, pETEC666

229    encodes a two-component sensor-kinase, herein designated *etcA* and *etcB* (***E. coli***

9

230   **t**wo-**c**omponent), and a three gene locus (herein designated *eor* for ***E. coli***

231   **o**xido**r**eductase) encoding a protein with homology to cytochrome b-type subunit

232   oxidoreductase protein (*eorA*), a protein with homology to an oxidoreductase

233   molybdopterin binding domain protein (*eorB*) and a periplasmic protein of unknown

234   function (*eorC*). In addition, ETEC H10407 pETEC58 encodes a putative

235   deoxycytidylate deaminase (pETEC58_0005).

236   As above, if a particular protein plays an important role in ETEC-mediated disease

237   then one would expect it to have a wide distribution amongst ETEC strains. To

238   determine whether the genes encoding the putative and known virulence factors of

239   the ETEC H10407 plasmids, which we identified above, were conserved amongst

240   ETEC strains we next examined their prevalence amongst the available sequenced

241   strains.  To aid this process we determined the sequence of the plasmids from ETEC

242   strain E1392/75.  *E. coli* E1392/75 possesses five plasmids three large conjugative

243   plasmids designated pETEC1018, pETEC746 and pETEC557 and two mobilizable

244   plasmids termed pETEC75 and pETEC62 (see Table 2 for general characteristics).

245   Included in the prevalence investigations were the ETEC strains E24377A and B7A

246   and the plasmid pCoo from ETEC strain C921b-1, all of which were sequenced in

247   other projects (28, 44). As the ETEC B7A genome is incomplete and no plasmids

248   were resolved, and pCoo is the only plasmid sequenced from ETEC C921b-1, we

249   can only confirm the presence of genes amongst the available DNA sequences and

250   not the absence of particular genes from these strains.  The distribution and location

251   of the cargo genes encoding known or putative virulence factors amongst the

252   sequenced ETEC plasmids is depicted in Fig. 3 and also listed in Table S4.

253   Comparative analyses revealed that, like ETEC H10407, the ETEC strains

254   E1392/75, B7A and E24377A possess the ST and LT enterotoxins (none were

255   identified for *E. coli* C921b-1, but previous analyses showed this strain to harbour LT

256   and ST) (54). The EtpABC two-partner secretion system was identified in ETEC

257   E1392/75 and E24377A; homologues may exist in ETEC strains B7A and C921b-1

258   but their existence or non existence in these strains could not be resolved due to the

259   lack of complete sequence data however other studies have not demonstrated a

260   universal association of the *etpABC* locus with ETEC strains (23). Unlike ETEC

261   strains H10407, E24377A and C921b-1, the autotransporter-encoding *eatA* gene

262   was not present on the ETEC E1392/75 plasmids.  A homologue annotated as EatA

263   is found in *E. coli* B7A, however further analyses of this protein reveal that it is more

264   closely related to SepA, a homologous SPATE protein from *Shigella flexneri* (1). No

265   equivalents of the ETEC H10407 *etcAB*, *eorA-C* or of the gene encoding the putative

266   deoxycytidylate deaminase, were detected in any of the other ETEC strains.

267   Like *E. coli* H10407, the ETEC strains E24377A, E1392/75 and C921b-1 encode

268   dispersin-like proteins previously designated CexE (43).  Further analyses reveal

269   that CexE is present in ETEC strains 27D and G427 (two CFA/I$^+$ strains) (43) and

270   ETEC O167:H5, a CS6 and CS5 encoding strain (9). For EAEC, dispersin is

271   secreted via the Aat Type I secretion system, associates non-covalently with the

272   extracellular face of the outer membrane preventing collapse of the AAF/II fimbriae

273   onto the bacterial cell surface by alteration of the surface charge and is required for

274   colonisation (31, 40, 52).  Analyses of the nucleotide sequences from ETEC strains

275   B7A, E24377A and E1392/75 reveal the presence of loci encoding Type I secretion

276   systems bearing striking homology to the Aat dispersin secretion system (Fig. 4).

277   The co-occurrence of *cexE* genes with *aat* loci suggests that the CexE proteins are

278   substrates for the Aat-like secretion systems of ETEC. Since, plasmid-borne fimbrial

279   loci are inextricably linked to ETEC-mediated disease (18), CexE may play a similar

11

280    role to dispersin by maintaining the CFs in a manner such that they can interact with

281    epithelial receptors. However, further studies are required to investigate the function

282    and distribution of CexE and to identify other relatives of this protein hitherto not

283    recognised.

284    As mentioned above, adherence via plasmid-encoded fimbrial systems is a crucial

285    step in ETEC pathogenesis (62).  *E. coli* H10407 pETEC948 possesses the CFA/I

286    chaperone-usher system (Fig. 3). ETEC E24377A possesses two-chaperone-usher

287    fimbrial systems located on pETEC_80 and pETEC_73 encoding the CS3 and CS1

288    fimbriae, respectively (44). Similarly, *E. coli* E1392/75 possesses the CS3- and CS1-

289    encoding loci on plasmids pETEC1018 and pETEC746 respectively, whereas pCoo

290    possesses the CS1 cluster, all of which have been described previously (28, 57, 58).

291    In addition, *E. coli* E1392/75 pETEC557 also encodes the CFA/III type IV fimbrium

292    (29). To determine whether fimbrial systems other than those mentioned above

293    might play a crucial role in ETEC pathogenesis we investigated conservation of

294    putative fimbrial loci amongst the available *E. coli* sequences. ETEC H10407

295    contains 12 additional loci predicted to encode fimbriae, all of which were

296    chromosomally located (Table S5). Four of these loci (*mat, sfm, ycb* and *yde*)

297    contain pseudogenes and were considered non-functional. We sought to establish if

298    *E. coli* H10407 harboured ETEC-specific fimbrial loci that might not be expressed by

299    commensal *E. coli, E. coli* K-12 or enteroaggregative *E. coli.*   The vast majority of

300    fimbrial operons identified are also located in commensal and laboratory strains with

301    notable exceptions.  The *yqi* and *stf-mrf* fimbrial loci are present in *E. coli* H10407

302    but contain pseudogenes in commensal or laboratory *E. coli.* However, an

303    apparently functional *yqi* operon is also present in enteroaggregative *E. coli* strain

304    042 and thus a functional *yqi* locus does not appear to be ETEC-specific. Indeed*,* the

12

305 *yqi* operon does not appear to be present in ETEC B7A (4). With regard to the *stf-mrf*

306 operon, the *mrfC* gene is a pseudogene in *E. coli* K-12 but not in ETEC H10407.

307 This six gene cluster (*smfA-mrfCD-stfEFG*) is present in ETEC E24377A and EAEC

308 042, though with some divergence in the *stf* genes.

309 Finally, the ETEC E1392/75 pETEC62 plasmid possesses CDSs encoding a type II

310 dihydropteroate synthase gene conferring sulfonamide resistance, and CDSs

311 encoding streptomycin phosphotransferase genes conferring streptomycin

312 resistance. This plasmid possesses 99% nucleotide identity with the ETEC E24377A

313 pETEC_6 plasmid and shares high levels of identity with plasmids from a variety of

314 *E. coli* sp including the *Shigella sonnei* pKKTET7 and the EPEC pE2348-2 plasmids

315 However, this plasmid has no homologue in ETEC H10407 and no detectable

316 homology amongst the ETEC B7A sequences suggesting it may not be widespread

317 amongst ETEC strains and thus is not essential for ETEC mediated diarrhea.

318 In conclusion, the putative and known virulence genes identified on the plasmids of

319 *E. coli* H10407 have a differential distribution amongst the sequenced ETEC strains.

320 In all cases the ETEC strains possess genes encoding the the ST and/or LT toxins

321 (*sta* and/or *eltAB, respectively*), a chaperone-usher fimbrial biogenesis locus (e.g.

322 the *cfa* locus) and components of an *aat-cexE* dispersin-like Type I secretion system.

323 Thus, despite the variation in individual plasmid gene content, comparison of the

324 entire plasmid complement of the sequenced ETEC strains suggests that there is a

325 conserved core of genes contained on the plasmids that are predicted to be involved

326 in virulence and may be essential for the establishment of ETEC-mediated disease.

327 **ETEC plasmids demonstrate a mosaic structure**

13

328　To determine whether the virulence factors identified above were encoded on a

329　specific plasmid, or repertoire of plasmids, we examined the nucleotide sequence

330　identity shared by the ETEC plasmids.  The nucleotide sequence of the conjugative

331　plasmids from each of the ETEC strains H10407, E1392/75 and E24377A were

332　concatenated and compared by BLASTn. The level of nucleotide sequence identity

333　between pCoo and the other ETEC plasmids was determined in a similar manner.

334　These comparisons revealed that while the plasmids all belong to a narrow subset of

335　incompatibility groups (see below), extensive rearrangements and recombination

336　events have occurred, resulting in individual plasmids that vary in their repertoires of

337　virulence genes (Fig. 3 and Table S4).　Such recombination is exemplified by

338　examining the distribution of the *eatA* gene. Thus, the *eatA* gene is not present in

339　ETEC strain 1392/75, in ETEC strain E24377A the *eatA* gene is located on

340　pETEC_74 and the *eltAB*, *aatPABC* and *etpABC* loci are located on pETEC_80. In

341　contrast, in ETEC strain H10407 the *eatA* gene is collocated with *etpABC* and

342　*aatPABC* on pETEC948 whereas the *eltAB* locus is located on pETEC666. The eatA

343　gene is present on ETEC C921b-1 pCoo along with *cooABCD* however in ETEC

344　strain E24377A *cooABCD* is located on a separate plasmid (pETEC_73) (Fig. 3 and

345　Table S4).  Other virulence-associated genes also display such differential

346　distribution (Table S4) suggesting that the extrachromosomal components of the

347　ETEC genome are in a state of flux (34, 44).　Notably, the plasmids contain an

348　extensive repertoire of IS elements and transposons (Table 2)(34); it is likely that the

349　mobility of these genetic elements, or the recombination between these elements,

350　gives rise to the observed mosaic structure of the ETEC plasmids.

351　Similar comparisons of the small mobilizable plasmids of the ETEC strains did not

352　demonstrate recombination between the mobilizsable plasmids. Furthermore, there

14

353    did not appear to be any significant exchange of genetic material between the

354    conjugative plasmids and the small mobilizable plasmids (data not shown).

355    **Plasmid stability and maintenance functions of the ETEC plasmids**

356    To determine whether the virulence factors described above were encoded on self

357    transmissible plasmids we examined the CDSs encoding the plasmid maintenance

358    and transfer functions of each ETEC plasmid.  A complete description of *E. coli*

359    H10407 pETEC666 has been published previously (41) and the complete repertoire

360    of genes for each ETEC plasmid are given in the EMBL databases (see Table 2 for

361    accession numbers) thus only the most salient features are described here. Plasmid

362    nomenclature utilises a system based on incompatibility groupings; plasmids of the

363    same incompatibility group should not co-exist within the same bacterial cell because

364    of the similarity in their replication systems. (34).  However, sequence analyses of

365    the CDSs encoding the plasmid replication functions of the repertoire of ETEC

366    plasmids revealed that the large conjugative-like plasmids of *E. coli* strains H10407,

367    E1392/75 and E24377A belong to a narrow subset of incompatibility groups and

368    possess multiple plasmids with the same replication mechanism (Fig. 3 and Table 2).

369    Thus, *E. coli* H10407 plasmids pETEC948 and pETEC666 belong to the RepFIIA

370    (IncFIIA) subset of incompatibility groupings and have RepA1 proteins which share

371    94% identity (95% similarity), whereas *E. coli* E1392/75 plasmids pETEC746 and

372    pETEC557 harbour RepI1 (IncI1) replication functions (*E. coli* E1392/75 pETEC557

373    is an apparent cointegrate of a RepF1B and RepI1 pasmids; such cointegration has

374    previously been noted for *E. coli* C921b-1, where pCoo represents a co-integrate

375    between a RepFIIA and a RepI1 plasmid (28)) with the corresponding RepZ proteins

376    sharing 94% identity (95% similarity). Similarly, the previously described ETEC strain

377    E24377A (44) possesses three plasmids with RepFIIA functions.  The basis for these

15

378   anti-dogmatic observations is not understood and requires further in depth

379   investigation.

380   Analyses of the nucleotide sequences of the repertoire of large conjugative-like

381   plasmids revealed that they possessed a number of plasmids stability systems

382   including post-segregation killing systems and active partitioning systems.   The

383   distribution of these systems amongst the plasmids sequenced in this study is given

384   in Table 2.   These stability systems have been described previously (see reviews

385   references (25, 56).

386   Previous studies have noted that the large plasmids encoding the toxins of ETEC are

387   in some cases self transmissible and in other cases not transmissible (27).   To

388   investigate whether the plasmids sequenced in this study possessed transmissibility

389   functions we examined the transfer regions of the conjugative-like plasmids.   As

390   noted previously, *E. coli* H10407 pETEC666 has a transfer region which is

391   interrupted by several IScE8 elements severely diminishing the ability of this system

392   to function efficiently (41).   In contrast, *E. coli* H10407 pETEC948 only possesses

393   remnants of the conjugation apparatus and is presumably not self transmissible.   In

394   addition, the *E. coli* E1392/75 pETEC1018 plasmid also contains an incomplete

395   conjugation apparatus which is presumed to be ineffective at promoting conjugation,

396   however *E. coli* E1392/75 pETEC746 possess an intact conjugation system that is

397   100% identical to the region encoding the functional R64-like conjugative pilus of

398   pCoo of *E. coli* C921b-1 and thus it is presumed to be functional. *E. coli* E1392/75

399   pETEC557 lacks CDSs encoding the R64 conjugative pilus and possesses remnants

400   of an F-like conjugation system.

16

401 ETEC strains H10407, E1392/75, and E24377A all contain similar small mobilizable

402 plasmids (pETEC52, pETEC75, and pETEC_5, respectively) with *mob* and *rep*

403 regions displaying 100% identity. The *E. coli* E1392/75 pETEC75 plasmid contains

404 an IS100 element not present in the other two plasmids. The distribution of these

405 plasmid types among the sequenced ETEC suggests that they might be common

406 components of ETEC genomes.   This plasmid type has been found in a number of

407 other *E. coli* strains and has been shown to increase the fitness of certain *E. coli* host

408 strains (16).   Therefore, multiple selective advantages might be conferred on the

409 ETEC strains possessing these small plasmids.  The *rep* and *mob* regions (3058 bp)

410 of ETEC H10407 pETEC58 plasmid, which encodes the putative deoxycytidylate

411 deaminase, demonstrates 81% identity with plasmid pHW66 from Rahnella sp.

412 WMR66; the putative deoxycytidylate deaminase is lacking from pHW66.  In contrast

413 to the other ETEC plasmids, there are no plasmids homologous to ETEC H10407

414 pETEC58 amongst the other genome sequenced ETEC isolates.

415

416 **The *E. coli* E1392/75 pETEC746 plasmid contains a pilin shufflon.**

417 As mentioned above, ETEC E1392/75 pETEC746 contains regions homologous to

418 the *S. enterica* Typhimurium RepI1 plasmid R64 that are also present in *E. coli*

419 C921b-1 pCoo and have been shown to be functional in that system (28).  During the

420 finishing of the ETEC genome, dideoxy sequencing of the region from 56,253 bp to

421 59,961 bp of pETEC746 from *E. coli* E1392/75 identified a nucleotide region

422 undergoing dynamic alteration.  The region of DNA consisted of a shufflon similar to

423 that of R64 (36).  PilV is a component of a conjugative pilus that expresses different

424 tips involved with attachment to cells.  The tips are regulated *via* a DNA shufflon

425 mechanism involving recombination at particular repeating sites.  Recombination is

17

426    mediated by the *rci* recombinase linked to this region.  Alternative tip adhesins are

427    involved in attachment to different strains and species and has been elucidated

428    experimentally in *S.* Typhimurium (36).  Evidence that the shufflon is functional in the

429    *E.coli* E1392/75 plasmid pETEC746 is provided in the sequences from a small insert

430    library.  Within the sequences are examples of *pilV* with alternative C-terminal tips,

431    implying that the plasmids sequenced represented a population in genetic flux.

432    There is direct evidence for sequences of *pilV* with tips *V1*, *V3* and *V4* (Fig. 5).

433    There are also regions of DNA sequence equivalent to tips *shuC1, shuC'* and *shuC2*

434    from *S.* Typhimurium.  However, these were only present in a small subpopulation of

435    pETEC746 plasmids and have been discounted from the complete finished

436    sequence.

437    **Conclusion**

438    This study provides a genomic context for the vast amount of experimental and

439    epidemiological data published thus far and provides a template for future diagnostic

440    and intervention strategies.  Evidence presented here suggests the prototypical

441    ETEC isolate *E. coli* H10407 was a commensal isolate that acquired a number of

442    plasmids containing a limited repertoire of virulence genes and thereby gained the

443    ability to cause disease.  Furthermore, comparisons of the genetic content of *E. coli*

444    H10407 with other ETEC strains reveals only a limited number of conserved genes

445    suggesting that to become pathogenic *E. coli* need only acquire (i) toxins (ST, LT or

446    both) to elicit net secretion from enterocytes, (ii) a fimbrial system that mediates

447    attachment to the intestinal epithelium e.g. CFA/I, and (iii) a novel Type I secretion

448    system the substrate of which (CexE) maintains the fimbriae in the correct physical

449    organisation.  This data suggests ETEC vaccine strategies should focus on these

450    plasmid-encoded virulence factors. However, given the relative plasticity of the *E.*

451     *coli* genome molecular epidemiological studies are essential to determine whether

452     these factors are widely distributed amongst ETEC strains from geographically

453     diverse locations.

454     **ACKNOWLEDGEMENTS**

457

## REFERENCES

459 1. Benjelloun-Touimi, Z., P. J. Sansonetti, and C. Parsot. 1995. SepA, the major extracellular
460 protein of Shigella flexneri: autonomous secretion and involvement in tissue invasion. Mol
461 Microbiol 17:123-35.
462 2. Brown, N. L., S. R. Barrett, J. Camakaris, B. T. Lee, and D. A. Rouch. 1995. Molecular
463 genetics and transport analysis of the copper-resistance determinant (pco) from
464 Escherichia coli plasmid pRJ1004. Mol Microbiol 17:1153-66.
465 3. Chang, D. E., D. J. Smalley, D. L. Tucker, M. P. Leatham, W. E. Norris, S. J. Stevenson, A. B.
466 Anderson, J. E. Grissom, D. C. Laux, P. S. Cohen, and T. Conway. 2004. Carbon nutrition of
467 Escherichia coli in the mouse intestine. Proc Natl Acad Sci U S A 101:7427-32.
468 4. Chaudhuri, R. R., M. Sebaihia, J. L. Hobman, M. A. Webber, D. L. Leyton, M. D. Goldberg, A.
469 F. Cunningham, A. Scott-Tucker, P. R. Ferguson, C. M. Thomas, G. Frankel, C. M. Tang, E. G.
470 Dudley, I. S. Roberts, D. A. Rasko, M. J. Pallen, J. Parkhill, J. P. Nataro, N. R. Thomson, and
471 I. R. Henderson. 2010. Complete genome sequence and comparative metabolic profiling of
472 the prototypical enteroaggregative Escherichia coli strain 042. PLoS One 5:e8801.
473 5. Chen, Q., S. J. Savarino, and M. M. Venkatesan. 2006. Subtractive hybridization and optical
474 mapping of the enterotoxigenic Escherichia coli H10407 chromosome: isolation of unique
475 sequences and demonstration of significant similarity to the chromosome of E. coli K-12.
476 Microbiology 152:1041-54.
477 6. Coster, T. S., M. K. Wolf, E. R. Hall, F. J. Cassels, D. N. Taylor, C. T. Liu, F. C. Trespalacios, A.
478 DeLorimier, D. R. Angleberger, and C. E. McQueen. 2007. Immune response, ciprofloxacin
479 activity, and gender differences after human experimental challenge by two strains of
480 enterotoxigenic Escherichia coli. Infect Immun 75:252-9.
481 7. Cravioto, A. 1980. Ph.D. Thesis. University of London, London.
482 8. Dallas, W. S. 1990. The heat-stable toxin I gene from Escherichia coli 18D. J Bacteriol
483 172:5490-3.
484 9. de Haan, L. A., G. A. Willshaw, B. A. van der Zeijst, and W. Gaastra. 1991. The nucleotide
485 sequence of a regulatory gene present on a plasmid in an enterotoxigenic Escherichia coli
486 strain of serotype O167:H5. FEMS Microbiol Lett 67:341-6.
487 10. Delcher, A. L., K. A. Bratke, E. C. Powers, and S. L. Salzberg. 2007. Identifying bacterial
488 genes and endosymbiont DNA with Glimmer. Bioinformatics 23:673-9.
489 11. Delcher, A. L., D. Harmon, S. Kasif, O. White, and S. L. Salzberg. 1999. Improved microbial
490 gene identification with GLIMMER. Nucleic Acids Res 27:4636-41.
491 12. Duriez, P., O. Clermont, S. Bonacorsi, E. Bingen, A. Chaventre, J. Elion, B. Picard, and E.
492 Denamur. 2001. Commensal Escherichia coli isolates are phylogenetically distributed
493 among geographically distinct human populations. Microbiology 147:1671-6.
494 13. Elsinghorst, E. A., and D. J. Kopecko. 1992. Molecular cloning of epithelial cell invasion
495 determinants from enterotoxigenic Escherichia coli. Infect Immun 60:2409-17.
496 14. Elsinghorst, E. A., and J. A. Weitz. 1994. Epithelial cell invasion and adherence directed by
497 the enterotoxigenic Escherichia coli tib locus is associated with a 104-kilodalton outer
498 membrane protein. Infect Immun 62:3463-71.
499 15. Emanuelsson, O., S. Brunak, G. von Heijne, and H. Nielsen. 2007. Locating proteins in the
500 cell using TargetP, SignalP and related tools. Nat Protoc 2:953-71.
501 16. Enne, V. I., P. M. Bennett, D. M. Livermore, and L. M. Hall. 2004. Enhancement of host
502 fitness by the sul2-coding plasmid p9123 in the absence of selective pressure. J Antimicrob
503 Chemother 53:958-63.
504 17. Evans, D. G., T. K. Satterwhite, D. J. Evans, Jr., and H. L. DuPont. 1978. Differences in
505 serological responses and excretion patterns of volunteers challenged with

506       enterotoxigenic Escherichia coli with and without the colonization factor antigen. Infect
507       Immun 19:883-8.

508 **18.** Evans, D. G., R. P. Silver, D. J. Evans, Jr., D. G. Chase, and S. L. Gorbach. 1975. Plasmid-
509       controlled colonization factor associated with virulence in Esherichia coli enterotoxigenic
510       for humans. Infect Immun 12:656-67.

511 **19.** Evans, D. J., Jr., and D. G. Evans. 1973. Three characteristics associated with
512       enterotoxigenic Escherichia coli isolated from man. Infect Immun 8:322-8.

513 **20.** Fleckenstein, J. M., P. R. Hardwidge, G. P. Munson, D. A. Rasko, H. Sommerfelt, and H.
514       Steinsland. Molecular mechanisms of enterotoxigenic Escherichia coli infection. Microbes
515       Infect 12:89-98.

516 **21.** Fleckenstein, J. M., D. J. Kopecko, R. L. Warren, and E. A. Elsinghorst. 1996. Molecular
517       characterization of the tia invasion locus from enterotoxigenic Escherichia coli. Infect
518       Immun 64:2256-65.

519 **22.** Fleckenstein, J. M., L. E. Lindler, E. A. Elsinghorst, and J. B. Dale. 2000. Identification of a
520       gene within a pathogenicity island of enterotoxigenic Escherichia coli H10407 required for
521       maximal secretion of the heat-labile enterotoxin. Infect Immun 68:2766-74.

522 **23.** Fleckenstein, J. M., K. Roy, J. F. Fischer, and M. Burkitt. 2006. Identification of a two-
523       partner secretion locus of enterotoxigenic Escherichia coli. Infect Immun 74:2245-58.

524 **24.** Francia, M. V., A. Varsaki, M. P. Garcillan-Barcia, A. Latorre, C. Drainas, and F. de la Cruz.
525       2004. A classification scheme for mobilization regions of bacterial plasmids. FEMS
526       Microbiol Rev 28:79-100.

527 **25.** Friehs, K. 2004. Plasmid copy number and plasmid stability. Adv Biochem Eng Biotechnol
528       86:47-82.

529 **26.** Frishman, D., A. Mironov, H. W. Mewes, and M. Gelfand. 1998. Combining diverse
530       evidence for gene recognition in completely sequenced bacterial genomes. Nucleic Acids
531       Res 26:2941-7.

532 **27.** Froehlich, B., E. Holtzapple, T. D. Read, and J. R. Scott. 2004. Horizontal transfer of CS1
533       pilin genes of enterotoxigenic Escherichia coli. J Bacteriol 186:3230-7.

534 **28.** Froehlich, B., J. Parkhill, M. Sanders, M. A. Quail, and J. R. Scott. 2005. The pCoo plasmid of
535       enterotoxigenic Escherichia coli is a mosaic cointegrate. J Bacteriol 187:6509-16.

536 **29.** Gomez-Duarte, O. G., S. Chattopadhyay, S. J. Weissman, J. A. Giron, J. B. Kaper, and E. V.
537       Sokurenko. 2007. Genetic diversity of the gene cluster encoding longus, a type IV pilus of
538       enterotoxigenic Escherichia coli. J Bacteriol 189:9145-9.

539 **30.** Gupta, A., K. Matsui, J. F. Lo, and S. Silver. 1999. Molecular basis for resistance to silver
540       cations in Salmonella. Nat Med 5:183-8.

541 **31.** Harrington, S. M., J. Sheikh, I. R. Henderson, F. Ruiz-Perez, P. S. Cohen, and J. P. Nataro.
542       2009. The Pic protease of enteroaggregative Escherichia coli promotes intestinal
543       colonization and growth in the presence of mucin. Infect Immun 77:2465-73.

544 **32.** Henderson, I. R., F. Navarro-Garcia, and J. P. Nataro. 1998. The great escape: structure and
545       function of the autotransporter proteins. Trends Microbiol 6:370-8.

546 **33.** Iguchi, A., N. R. Thomson, Y. Ogura, D. Saunders, T. Ooka, I. R. Henderson, D. Harris, M.
547       Asadulghani, K. Kurokawa, P. Dean, B. Kenny, M. A. Quail, S. Thurston, G. Dougan, T.
548       Hayashi, J. Parkhill, and G. Frankel. 2009. Complete genome sequence and comparative
549       genome analysis of enteropathogenic Escherichia coli O127:H6 strain E2348/69. J Bacteriol
550       191:347-54.

551 **34.** Johnson, T. J., and L. K. Nolan. 2009. Pathogenomics of the virulence plasmids of
552       Escherichia coli. Microbiol Mol Biol Rev 73:750-74.

553 **35.** Kaper, J. B., J. P. Nataro, and H. L. Mobley. 2004. Pathogenic Escherichia coli. Nat Rev
554       Microbiol 2:123-40.

555 **36.** Komano, T., S. R. Kim, and T. Yoshida. 1995. Mating variation by DNA inversions of
556       shufflon in plasmid R64. Adv Biophys 31:181-93.

557   **37.**   **Krogh, A., B. Larsson, G. von Heijne, and E. L. Sonnhammer. 2001. Predicting**
558          **transmembrane protein topology with a hidden Markov model: application to complete**
559          **genomes. J Mol Biol 305:567-80.**

560   **38.**   **McHugh, G. L., R. C. Moellering, C. C. Hopkins, and M. N. Swartz. 1975. Salmonella**
561          **typhimurium resistant to silver nitrate, chloramphenicol, and ampicillin. Lancet 1:235-40.**

562   **39.**   **Nataro, J. P., and J. B. Kaper. 1998. Diarrheagenic Escherichia coli. Clin Microbiol Rev**
563          **11:142-201.**

564   **40.**   **Nishi, J., J. Sheikh, K. Mizuguchi, B. Luisi, V. Burland, A. Boutin, D. J. Rose, F. R. Blattner,**
565          **and J. P. Nataro. 2003. The export of coat protein from enteroaggregative Escherichia coli**
566          **by a specific ATP-binding cassette transporter system. J Biol Chem 278:45680-9.**

567   **41.**   **Ochi, S., T. Shimizu, K. Ohtani, Y. Ichinose, H. Arimitsu, K. Tsukamoto, M. Kato, and T. Tsuji.**
568          **2009. Nucleotide sequence analysis of the enterotoxigenic Escherichia coli Ent plasmid.**
569          **DNA Res 16:299-309.**

570   **42.**   **Patel, S. K., J. Dotson, K. P. Allen, and J. M. Fleckenstein. 2004. Identification and**
571          **molecular characterization of EatA, an autotransporter protein of enterotoxigenic**
572          **Escherichia coli. Infect Immun 72:1786-94.**

573   **43.**   **Pilonieta, M. C., M. D. Bodero, and G. P. Munson. 2007. CfaD-dependent expression of a**
574          **novel extracytoplasmic protein from enterotoxigenic Escherichia coli. J Bacteriol 189:5060-**
575          **7.**

576   **44.**   **Rasko, D. A., M. J. Rosovitz, G. S. Myers, E. F. Mongodin, W. F. Fricke, P. Gajer, J. Crabtree,**
577          **M. Sebaihia, N. R. Thomson, R. Chaudhuri, I. R. Henderson, V. Sperandio, and J. Ravel.**
578          **2008. The pangenome structure of Escherichia coli: comparative genomic analysis of E. coli**
579          **commensal and pathogenic isolates. J Bacteriol 190:6881-93.**

580   **45.**   **Ren, C. P., R. R. Chaudhuri, A. Fivian, C. M. Bailey, M. Antonio, W. M. Barnes, and M. J.**
581          **Pallen. 2004. The ETT2 gene cluster, encoding a second type III secretion system from**
582          **Escherichia coli, is present in the majority of strains but has undergone widespread**
583          **mutational attrition. J Bacteriol 186:3547-60.**

584   **46.**   **Roy, K., G. M. Hilliard, D. J. Hamilton, J. Luo, M. M. Ostmann, and J. M. Fleckenstein. 2009.**
585          **Enterotoxigenic Escherichia coli EtpA mediates adhesion between flagella and host cells.**
586          **Nature 457:594-8.**

587   **47.**   **Rutherford, K., J. Parkhill, J. Crook, T. Horsnell, P. Rice, M. A. Rajandream, and B. Barrell.**
588          **2000. Artemis: sequence visualization and annotation. Bioinformatics 16:944-5.**

589   **48.**   **Sammut, S. J., R. D. Finn, and A. Bateman. 2008. Pfam 10 years on: 10,000 families and still**
590          **growing. Brief Bioinform 9:210-9.**

591   **49.**   **Schubert, S., A. Rakin, and J. Heesemann. 2004. The Yersinia high-pathogenicity island**
592          **(HPI): evolutionary and functional aspects. Int J Med Microbiol 294:83-94.**

593   **50.**   **Scotland, S. M., N. P. Day, and B. Rowe. 1983. Acquisition and maintenance of enterotoxin**
594          **plasmids in wild-type strains of Escherichia coli. J Gen Microbiol 129:3111-20.**

595   **51.**   **Selvaratnam, S., and M. A. Gealt. 1993. Transcription of ColE1Ap mbeC induced by**
596          **conjugative plasmids from twelve different incompatibility groups. J Bacteriol 175:6982-7.**

597   **52.**   **Sheikh, J., J. R. Czeczulin, S. Harrington, S. Hicks, I. R. Henderson, C. Le Bouguenec, P.**
598          **Gounon, A. Phillips, and J. P. Nataro. 2002. A novel dispersin protein in enteroaggregative**
599          **Escherichia coli. J Clin Invest 110:1329-37.**

600   **53.**   **Silver, S., A. Gupta, K. Matsui, and J. F. Lo. 1999. Resistance to ag(i) cations in bacteria:**
601          **environments, genes and proteins. Met Based Drugs 6:315-20.**

602   **54.**   **Smyth, C. J. 1982. Two mannose-resistant haemagglutinins on enterotoxigenic Escherichia**
603          **coli of serotype O6:K15:H16 or H-isolated from travellers' and infantile diarrhoea. J Gen**
604          **Microbiol 128:2081-96.**

605   **55.**   **Steffen, R., F. Castelli, H. Dieter Nothdurft, L. Rombo, and N. Jane Zuckerman. 2005.**
606          **Vaccination against enterotoxigenic Escherichia coli, a cause of travelers' diarrhea. J Travel**
607          **Med 12:102-7.**

608     **56.**     **Summers, D. K., and D. J. Sherratt. 1985. Bacterial plasmid stability. Bioessays 2:209-211.**

609     **57.**     **Svennerholm, A. M., and C. Ahren. 1982. Serological subtypes of Escherichia coli**
610          **colonization factor antigen II. Eur J Clin Microbiol 1:107-11.**

611     **58.**     **Tacket, C. O., R. H. Reid, E. C. Boedeker, G. Losonsky, J. P. Nataro, H. Bhagat, and R.**
612          **Edelman. 1994. Enteral immunization and challenge of volunteers given enterotoxigenic E.**
613          **coli CFA/II encapsulated in biodegradable microspheres. Vaccine 12:1270-4.**

614     **59.**     **Tenaillon, O., D. Skurnik, B. Picard, and E. Denamur. The population genetics of**
615          **commensal Escherichia coli. Nat Rev Microbiol 8:207-17.**

616     **60.**     **Turner, A. K., T. D. Terry, D. A. Sack, P. Londono-Arcila, and M. J. Darsley. 2001.**
617          **Construction and characterization of genetically defined aro omp mutants of**
618          **enterotoxigenic Escherichia coli and preliminary studies of safety and immunogenicity in**
619          **humans. Infect Immun 69:4969-79.**

620     **61.**     **Turner, S. M., R. R. Chaudhuri, Z. D. Jiang, H. DuPont, C. Gyles, C. W. Penn, M. J. Pallen, and**
621          **I. R. Henderson. 2006. Phylogenetic comparisons reveal multiple acquisitions of the toxin**
622          **genes by enterotoxigenic Escherichia coli strains of different evolutionary lineages. J Clin**
623          **Microbiol 44:4528-36.**

624     **62.**     **Turner, S. M., A. Scott-Tucker, L. M. Cooper, and I. R. Henderson. 2006. Weapons of mass**
625          **destruction: virulence factors of the global killer enterotoxigenic Escherichia coli. FEMS**
626          **Microbiol Lett 263:10-20.**

627     **63.**     **van der Woude, M. W., and I. R. Henderson. 2008. Regulation and Function of Ag43 (Flu).**
628          **Annu Rev Microbiol 62:153-169.**

629     **64.**     **Wenneras, C., and V. Erling. 2004. Prevalence of enterotoxigenic Escherichia coli-**
630          **associated diarrhoea and carrier state in the developing world. J Health Popul Nutr 22:370-**
631          **82.**

632     **65.**     **Yamamoto, T., T. Tamura, and T. Yokota. 1984. Primary structure of heat-labile**
633          **enterotoxin produced by Escherichia coli pathogenic for humans. J Biol Chem 259:5037-44.**

634     **66.**     **Yamamoto, T., and T. Yokota. 1980. Cloning of deoxyribonucleic acid regions encoding a**
635          **heat-labile and heat-stable enterotoxin originating from an enterotoxigenic Escherichia**
636          **coli strain of human origin. J Bacteriol 143:652-60.**

637

638

1 **Table 1.** General characteristics of three sequenced *E. coli* chromosomes

| Strain | H10407 | K-12 | HS |
|---|---|---|---|
| Etiology | Pathogen | Lab strain | Commensal |
| Length (bp) | 5153435 | 4643538 | 4686137 |
| GC content | 50.8% | 50.8% | 50.8% |
| Total CDS | 4746 | 4384 | 4200 |
| tRNA | 87 | 86 | 86 |
| rRNA | 7 | 7 | 7 |

2

Table 2. General characteristics of the plasmids from ETEC strains H10407 and E1392/75

| Strain | E. coli H10407 | | | | E. coli E1392/75 | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Plasmid | pETEC948 | pETEC666 | pETEC58 | pETEC52 | pETEC1018 | pETEC746 | pETEC557 | pETEC75 | pETEC62 |
| Accession No. | FN649418 | FN649417 | FN649416 | FN649415 | FN822745 | FN822748 | FN822746 | FN822749 | FN822747 |
| Size (bp) | 94797 | 66681 | 5800 | 5175 | 101857 | 74575 | 55709 | 7497 | 6222 |
| CDS | 115 | 88 | 7 | 6 | 165 | 117 | 73 | 9 | 13 |
| rep | RepFIIA | RepFIIA | ColE2 | ColE1 | RepFIIA | RepI1 | RepFIB/RepI1 | ColE1 | ND[a] |
| Stability genes | StbAB, PsiAB, SopAB, YacAB, RelE | StbAB PsiAB Mok/Hok | | | StbAB PsiAB CcdAB | StbAB NikAB | SopAB PsiAB | | |
| Insertion elements | IS1, IS2, IS3, IS66, IS91, IS100, IS629, IS911, IS1414, ISEc10, ISEc12, ISSfl4, Tn3 | IS1, IS21, IS66, IS600, IS1294, ISEc8 | | | IS1, IS2, IS3, IS21, IS30, IS66, IS91, IS100, IS629, IS630, IS639, IS911, IS1414, ISShdy1 | IS2, IS100, IS186, IS1328 | IS1, IS30, IS66, IS100, IS911, ISShdy1 | IS100 | ISCR2 |

*ND: not determined. pETEC62 has a gene conserved amongst many small plasmids which is annotated as a "probable replication initiation protein" but no experimental evidence exists for this function.

1    **Figure 1.** Circular representation of the *E. coli* H10407 chromosome. From the

2    outside in the outer circle 1 marks the position of regions of difference (mentioned in

3    the text) including prophage (light pink) as well as regions differentially present in

4    other *E. coli* strains: blue (See table S1). Circle 2 shows the size in bps. Circles 3

5    and 4 show the position of CDSs transcribed in a clockwise and anticlockwise

6    direction, respectively (for colour codes see below). Genes in circles 3 and 4 are

7    colour coded according to the function of their gene products: dark green=membrane

8    or surface structures, yellow=central or intermediary metabolism, cyan=degradation

9    of macromolecules, red=information transfer/cell division, cerise =degradation of

10    small molecules, pale blue =regulators, Salmon pink=pathogenicity or adaptation,

11    black=energy metabolism, orange=conserved hypothetical, pale green=unknown,

12    brown=pseudogenes. Circles 5 & 6 and 9 &10 show the position of *E. coli* H10407

13    genes which have orthologues (by reciprocal FASTA analysis) in *E. coli*: K-12

14    MG1655 (blue) or *E. coli* 042 (green), respectively. Circles 7 & 8 and 11 & 12 show

15    the position of genes unique to *E. coli* H10407 compared to *E. coli* K-12 MG1655

16    (red) or *E. coli* 042 (grey), respectively. Circle 13 shows a plot of G+C content (in a

17    10 Kb window). Circle 14 shows a plot of GC skew ([G-C]/[G+C]; in a 10 Kb window).

**Figure 2.** Comparison of the genetic content of *E. coli* H10407 chromosome with the chromosomes of other sequenced strains of *E. coli*. (A) Comparison of *E. coli* H10407 with the three non-pathogenic *E. coli* strains HS, C and K-12 reveals the four strains share a large proportion of common genes. Only 599 *E. coli* H10407 specific genes were identified. The *E. coli* H10407 specific CDS are not thought to be associated with virulence (see text for details). (B) Comparison of *E. coli* H10407 with the genome sequenced ETEC strains E24377A and B7A. The four strains possess 3553 genes in common however the ETEC strains share only 188 genes not present in the commensal strain *E. coli* HS. However, these latter genes are not unique to ETEC and are widely distributed amongst *E. coli* and are largely present among non-pathogenic strains of *E. coli* such as *E. coli* K-12

**Figure 3.** Nucleotide sequence comparison of large conjugative-like plasmids from ETEC strains. Plasmid sequences from each strain were concatenated and compared using BLASTn. BLAST matches longer than 250 bp are shown as grey blocks in a comparison between plasmids from E24377A (pETEC_80, pETEC_74, pETEC_73 and pETEC_35), H10407 (pETEC948 and pETEC666), E1392/75 (pETEC1018, pETEC746 and pETEC557) and C921b-1 (pCoo). Shading of the grey blocks is proportional to the BLAST match (minimum = 80% nucleotide identity, maximum = 100% nucleotide identity). Each plasmid is denoted as a linear black line, the identity of each plasmid is noted above the line and the source ETEC strain from which the plasmids are derived is given on the left side of the figure. Coding sequences are depicted by arrows and are coloured according to known or predicted function: blue, virulence-related; red, plasmid-related protein; green, outer membrane-related (includes conjugal transfer loci); pink, transposase/insertion element-related; light blue, regulatory protein; orange, conserved hypothetical protein; uncoloured, hypothetical protein. The position of genes encoding known or predicted virulence-related proteins is denoted by white boxes harbouring the gene names. In addition, the locus encoding the R64 conjugative pilus and the variant PilV tips is also depicted. The putative origin of replication associated with each of the plasmids is highlighted within yellow shaded boxes. The chimeric nature of the plasmids is clearly visible with recombination between plasmids a frequent occurrence. The unlabeled figure was prepared using a custom script (Sullivan MJ and SA Beatson, unpublished).

1

**Figure 4.** Comparison of the EAEC *aat-aap* locus with the *aat-cexE* loci of ETEC strains.

(A) The genetic organisation of the *aat* and *cexE* loci is depicted. The level of amino acid

identity for each component of the *aat-cexE* system is shown; figures represent comparison

with the *E. coli* H10407 orthologues.   Orthologues are coloured coded for ease of

identification. Genes which are not juxtaposed are depicted with a blue line separating them.

(B) Amino acid sequence alignment of ETEC CexE with the EAEC 042 dispersin.  All three

proteins possess a signal sequence which is cleaved after the amino acid at position 21 in

the  alignment.    There  is  limited  conservation  in  the  sequences  however  two  cysteine

residues which are disulphide bonded in dispersin are conserved.  Based on the structure of

dispersin, the remainder of the conserved residues appear to represent hydrophobic core

residues required for structural integrity of the molecule.

**Figure 5.** Arrangement of *pilV* shufflon region of *E. coli* E1392/75 pETEC746. Annotation of *pilV* region shown using the Artemis sequence viewer (1). Sequence blocks encoding C-terminal fragments of PilV are found in both orientations between *pilV* and the *rci* recombinase. Identical 13 bp repeats (gtgccaatccggt) are shown as miscellaneous features and mark the predicted sites of recombination between the C-terminal fragments and the *pilV* gene.

ETP507

ETP33

ETP468

ETP86

ETP128

E. coli H10407

ETP216

ETP295    ETP284

23  24  5000001

22

20

4000001

19

18

17

3000001

0

1

2

4

1000001

8

2000001

9

11

14  13  12

(A)



**MG1655**

206

65

**H10407**

599

102

41

22

3658

49

41

58

177

87

226

**C ATCC 8739**

**HS**

(B)



**H10407**

510

72

**HS**

215

192

25

73

3553

27

188

111

384

178

330

**E24377A**

**B7A**

(A)



EAEC 042    33%    27%    31%    50%    33%    18%

CexE

H10407    D    P    A    B    C

E1392/75    66%    63%    70%    77%    41%    96%

(B)

```
                    10         20         30         40         50         60
            ....|....|....|....|....|....|....|....|....|....|....|....|....|
H10407      MKKYILGVI--LAMGSLSAIAGGGNSERPPSVAAGECVTFNSKLGEIGGYSWKYSNDACN
E1392/75    MKKIILALP--FLTSCFSAFAGGSGPEWQPQISPGQCIQY-TEIGETGGYKWHN-IDACN
042         MKKIKFVIFSGILGISLNAFAGGSGWN-ADNVDPSQCIKQ-----SGVQYTYNSGVSVCM
            ***   : :     :  .:.*:***.. :   .: ..:*:       *.::    ..*
                    70         80         90        100        110        120
            ....|....|....|....|....|....|....|....|....|....|....|....|..
H10407      ETVAKGYAIGVAMHRTVNYEGGYSIQSSGIVKPGSDFIMKGGKTYKGHKKVSAGGDTPYWYK
E1392/75    EVVHRGYASGAFVSGKVVYEGGETIEYTGIVKPDAPYTIQAPSTHNGKKKVGHGGAYTYWAR
042         QGLNEGKVRGVSVSGVFYYNDGTTSNFKGVVTPSTPVNTNQDINKTNKVGVQKYRALTEWVK
            : : .* . *.:    . *:.* : : .*:*.*.:    :    . ..:  *      . * :
```