

1 **Recent O-antigen diversification masks highly pathogenic STEC O104:H4**

2

3 Christina Lang¹, Angelika Fruth¹, Ian Winsten Campbell², Claire Jenkins³, Peyton Smith⁴,
4 Nancy Strockbine⁴, François-Xavier Weill⁵, Ulrich Nübel^{6, 7, 8}, Yonatan H. Grad⁹, Matthew K.
5 Waldor^{2,9,10}, and Antje Flieger^{1#}

6

7 ¹ Robert Koch-Institut, Division of Enteropathogenic Bacteria and *Legionella*, National
8 Reference Centre for *Salmonella* and other Enteric Bacterial Pathogens, Burgstr. 37, 38855
9 Wernigerode, Germany

10 ² Brigham & Women's Hospital, Division of Infectious Diseases, Harvard Medical School,
11 Department of Microbiology, Boston, MA 02115, United States of America

12 ³ Gastro and Food Safety (One Health) Division, Health Security Agency, London NW9 5HT,
13 United Kingdom

14 ⁴ Centers for Disease Control and Prevention, Division of Foodborne, Waterborne and
15 Environmental Diseases, National Center for Emerging and Zoonotic Infectious Diseases,
16 1600 Clifton Road, Atlanta, GA, United States of America

17 ⁵ Institut Pasteur, Université Paris Cité, Unité des bactéries pathogènes entériques, Paris, F-
18 75015, France

19 ⁶ Leibniz Institute DSMZ - German Collection of Microorganisms and Cell Cultures,
20 Braunschweig, Germany

21 ⁷ German Center for Infection Research (DZIF), Partner Site Braunschweig-Hannover,
22 Germany

23 ⁸ Braunschweig Integrated Center of Systems Biology (BRICS), Technical University,
24 Braunschweig, Germany

25 ⁹ Department of Immunology and Infectious Diseases, Harvard T.H. Chan School of Public
26 Health, Boston, MA 02115, United States of America

27 ¹⁰ Howard Hughes Medical Institute, Boston, MA 02115, United States of America

28 # corresponding author contact: email: fliegera@rki.de, Tel.: +49-30-18754-2522

29 **Abstract**

30 Background: Shiga toxin-producing *E. coli* (STEC) can give rise to a range of clinical outcomes
31 from diarrhea to the life-threatening systemic condition, hemolytic uremic syndrome (HUS). A
32 major outbreak of HUS occurred in 2011, and was caused by a rare serotype, STEC O104:H4.
33 Prior to 2011 and since the outbreak, STEC O104:H4 were rarely associated with human
34 infections.

35 Methods: From 2012 to 2020 intensified STEC surveillance was performed in Germany where
36 subtyping of ~8,000 clinical isolates by molecular methods including whole genome
37 sequencing was carried out. Virulence traits and phylogenetic context were investigated for a
38 subset of strains.

39 Results: A rare STEC serotype O181:H4 associated with HUS was identified, belonging to
40 sequence type (ST) 678, like the STEC O104:H4 outbreak strain. Virulence and genomic
41 comparisons revealed that the two strains are phylogenetically related and differ principally in
42 the gene cluster encoding their respective lipopolysaccharide O-antigens. In addition, five
43 other serotypes belonging to ST678 from human clinical infection were identified from diverse
44 locations worldwide.

45 Conclusion: Our data suggest the high virulence ensemble of STEC O104:H4 remains a global
46 threat, but that horizontal exchange of O-antigen gene clusters has cloaked the pathogen with
47 new O-antigens, confounding interpretation of their potential risk.

48

49

50 **Keywords**

51 Shiga toxin-producing *E. coli*, O104:H4, O-antigen diversification, phylogeny, risk profiling

52

53

54 STEC are food-borne pathogens responsible for a range of clinical syndromes from diarrhea
55 to the life-threatening systemic condition, HUS, a triad of thrombotic microangiopathy,
56 thrombocytopenia, and acute renal injury [1]. Historically, the classification of STEC strains
57 into different serotypes has proven invaluable for epidemiology and risk profiling [2]. In *E. coli*
58 the serotype is determined by a combination of O- and H-antigen types (see below) whereas
59 the O group solely denotes the O-antigen. Globally, strains of STEC serotype O157:H7 are
60 most frequently associated with HUS, but STEC belonging to O groups O26, O103, O111, and
61 O145, have also been regularly linked to HUS development [2, 3]. In addition, a very rare
62 serotype gave rise to a major HUS outbreak in early summer 2011, when an O104:H4 strain
63 caused more than 3,000 cases of diarrhea and 800 cases of HUS including 54 fatalities,
64 predominantly in Germany [4-6].

65

66 The O104:H4 outbreak strain encodes an exceptional set of virulence features [5-9]. Like other
67 HUS-associated strains, the outbreak strain produces Shiga toxin (Stx), specifically the Stx2a
68 variant. But unlike most *E. coli* strains causing HUS, this strain is an enteroaggregative *E. coli*
69 (EAEC) and lacks a prime virulence factor, the type III secretion system encoded on a
70 pathogenicity island designated as the locus of enterocyte effacement (LEE). The O104:H4
71 outbreak strain along with other O104:H4 strains, all belonging to multi-locus sequence type
72 (MLST) ST678 and some possessing *stx*, formed a distinct clade among EAEC [5, 7, 10]. The
73 outbreak strain harbors a plasmid (pAA) characteristic for EAEC encoding aggregative
74 adhesion fimbriae of type I (AAF/I) [5, 7]. Furthermore, the outbreak strain encodes the
75 virulence-linked serine-protease autotransporters (SPATE), SepA, SigA, and Pic [5, 10, 11].
76 Despite the extensive outbreak in May/June of 2011 and associated wide distribution of the
77 strain in affected regions, intense molecular surveillance only uncovered relatively few
78 O104:H4 cases in Germany after the outbreak dissipated by July 2011 [12].

79

80 In *E. coli*, serotypes are determined by the composition of the lipopolysaccharide (LPS) O-
81 antigen and the flagellar H-antigen, both important surface features of microorganisms that

82 shape pathogen host interactions [13, 14]. LPS forms a major structural component of the
83 Gram-negative cell's outer membrane and its most distal part is the O-antigen. The O-antigen
84 is subject to strong selection pressure and is one of the most variable components of the
85 bacterial cell [13]. Typically, in *E. coli* the O-antigen consists of chains of repeating
86 oligosaccharide subunits, usually composed of two to seven sugars often with additional
87 chemical modifications [15]; currently, 182 O groups and 53 flagellar antigen types have been
88 described by phenotypic identification [14, 16]. The genes encoding O-antigen biosynthesis
89 are organized in clusters that are flanked by colanic acid (*wca*) and histidine (*his*) biosynthesis
90 genes [15]. O-antigen biosynthesis gene clusters typically have a lower GC content (often
91 <40%) than that of the backbone of the *E. coli* chromosome, which is ~50% GC content [13,
92 17, 18]. These differences suggest that O-antigen biosynthesis gene clusters are exchanged
93 by lateral gene transfer and are under diversifying selection, and therefore a hot spot of
94 recombination [15, 19].

95

96 Although its wide distribution in the affected areas, the near disappearance of *E. coli* O104:H4
97 in Germany after the large outbreak in 2011 was unanticipated. Here, we show that the high
98 virulence ensemble of the O104:H4 outbreak strain remains a threat, but that O-antigen gene
99 exchange has cloaked the pathogen with several new O-antigens.

100

101 **Methods**

102

103 **Ethical statement**

104 Rabbit studies conducted according to protocols reviewed and approved by Brigham and
105 Women's Hospital Committee on Animals (IACUC protocol 2016N000334) and Animal Welfare
106 Assurance of Compliance (number A4752-01) in accordance with recommendations in the
107 Guide for the Care and Use of Laboratory Animals of the National Institutes of Health and the
108 Animal Welfare Act of the U.S. Department of Agriculture.

109

110 **Study strains**

111 In the context of intensified STEC surveillance in Germany, clinical isolates were collected at
112 the National Reference Center for *Salmonella* and other Bacterial Enteric Pathogens and
113 analyzed for serotype, *stx* and subtypes, *eaeA*, *hlyA* and *aatA* as described [20]. Genome
114 sequencing was carried out on a subset of strains listed together with open source and study
115 strain data in Table S1 and S2. Strains were grown on nutrient agar (Oxoid GmbH, Germany),
116 Luria Bertani (LB) broth, or enterohemolysin agar (Sifin GmbH, Germany).

117

118 **Whole genome sequencing (WGS)**

119 Long read whole genome sequencing of O181:H4 strain 17-07187 was performed by GATC
120 Biotech (Konstanz, Germany) using a PacBio RS II sequencer (Pacific Biosciences, USA) and
121 short read genome sequencing of strains 12-04810, 14-03615, 14-01288, 16-00596, 16-
122 01499, 16-05332, 17-01774, 17-00416, 17-07187 and 19-02696 (Table S1) on an Illumina
123 MiSeq and HiSeq 1500 benchtop sequencer. Polishing of the assembled genome and
124 plasmids was performed with Illumina short reads by Pilon (version 1.22) [21]. The sequences
125 were uploaded to NCBI project: PRJNA833419.

126

127 **Bioinformatics analyses**

128 *De novo* assembly of the PacBio sequence data (103-fold average coverage) was performed
129 by GATC utilizing HGAP3 (Pacific Biosciences, USA). Quality control and trimming of MiSeq
130 raw reads with subsequent detection of serotype and virulence genes was performed as
131 described [16]. Genomic comparisons were carried out using MAUVE (version: 1.1.1) and
132 MAFFT (version 1.3.7) as plugin in Geneious (version: 11.1.5; Biomatters Ltd) [22, 23]. Ridom
133 SeqSphere+ (version:7.2.0, Ridom GmbH, Germany) was used to determine MLST Warwick
134 sequence types and to create minimal spanning trees based on 2513 allele targets from the
135 *E. coli* cgMLST Enterobase with pairwise ignoring missing values from genome assemblies
136 [24]. Phage prediction was carried out by analysing the genome sequences with PHASTER
137 [25]. RAST was used for CDS annotation [26].

138

139 **SNP-based alignment and maximum likelihood based phylogenetic tree**

140 Mapping of sequencing reads, generation of consensus sequences, and alignment calculation
141 was performed using the BatchMap pipeline [27]. The genome sequence of FWSEC0009 was
142 used as reference and SNPs were filtered using Gubbins (version: 3.2.1) [28]. Alignment of
143 filtered SNPs by Gubbins was used to generate a maximum likelihood based phylogenetic tree
144 by PHYML 3.3.20180214 (Geneious plugin, substitution model: HKY85, 100 bootstraps)[29].

145

146 **Temporal signal and 'clocklikeness' of molecular phylogenies**

147 TempEst was used to analyze the RAxML tree generated by Gubbins in conjunction with
148 collection year data to validate the molecular-clock assumption [30]. Best-fitting root was
149 chosen for linear regression analyses.

150

151 **Cytotoxicity, adherence, and infection assays**

152 Viability of Vero cells after inoculation with diluted bacterial culture supernatants (1:200) was
153 examined using 3-[4,5-dimethylthiazole-2-yl]-2,5-diphenyltetrazolium bromide [5]. Bacteria
154 and Vero cells were prepared as described [20]. Adherence to HEp-2 cells was performed as
155 described [27]. For infant rabbit infection assays, litters of mixed gender 2-day-old New
156 Zealand White infant rabbits with the lactating doe were acquired from Charles River (Canada,
157 strain code 052). Infant rabbits were orogastrically inoculated on the day of arrival with 10⁹
158 CFU of Streptomycin-resistant strains O104:H4 C227-11 and O181:H4 17-07187 suspended
159 in 500µl 2.5% sodium bicarbonate (pH9) using a size 4 French catheter as described
160 previously, except that no ranitidine was administered [11]. Infant rabbits were monitored for
161 signs of illness and euthanized three days post infection. Tissue samples taken from the
162 stomach, small intestine, cecum, and colon were homogenized and CFU determined by serial
163 dilution, and plating on LB media containing 200 µg/ml streptomycin [11].

164

165

166 **Results**

167 **The HUS-associated STEC O181:H4 strain 17-07187 shares a close phylogenetic**
168 **relationship and similar virulence traits with the O104:H4 outbreak strain**

169 After the large STEC O104:H4 outbreak in 2011, the German National Reference Centre for
170 *Salmonella* and other Bacterial Enteric Pathogens intensified STEC surveillance and analyzed
171 ~ 8,000 clinical isolates primarily from diarrhea and HUS patients from 2012 to 2020. This
172 strain collection included a stool sample isolate (17-07187) from a 6-year-old girl who had
173 bloody diarrhea and HUS in December 2017. She had not travelled outside her home in
174 Northwest Germany before becoming ill. Serotyping and whole genome sequencing revealed
175 that the strain belonged to an unusual serotype, O181:H4, that had not been previously
176 associated with HUS. The strain had *stx2a* but lacked the type III secretion system encoded
177 on the LEE pathogenicity island (marker gene *eae*) found in STEC. Furthermore, the strain
178 encoded characteristic EAEC markers including *aatA*, *aggR*, AAF/I genes and autotransporter
179 protease genes *pic*, *sigA*, and *sepA* (Fig.1A).

180
181 MLST demonstrated that the 17-07187 isolate belonged to the same sequence type (ST678)
182 as the O104:H4 outbreak strain (Table S2) [5, 7, 9]. Core genome MLST (cgMLST) [24]
183 confirmed the close phylogenetic relationship of this isolate with a panel of O104:H4 strains
184 (Fig.1B). There were only 30-34 allelic differences between this O181:H4 isolate and the
185 O104:H4 outbreak strain from 2011 (e.g. strains FWSEC0009, C227-11, or 11-02027), and
186 O104:H4 clinical isolates from the Republic of Georgia in 2009 (2009EL-2071) and France in
187 2012 (Ec12-0465). In contrast, there was considerable phylogenetic distance to STEC
188 serotypes O181:H16 (ST6274), O181:H49 (ST173), O104:H21 (ST672), and O104:H7
189 (ST10075) isolated between 2012 and 2019 (allelic distance >1500) (Fig.1B). Comparison of
190 the virulence gene repertoire of the O181:H4 and the O104:H4 outbreak strain also strongly
191 suggested that they rely on very similar virulence mechanisms (Fig.1A). Indeed, both strains
192 had comparable Stx-related cytotoxicity (Fig.1C), exhibited a characteristic enteroaggregative
193 adherence pattern (Fig.1D), and colonized intestinal tissues, particularly the cecum and colon,

194 similarly during *in vivo* infant rabbit infections (Fig.1E). Thus, STEC O181:H4 and O104:H4
195 isolates share marked genomic similarity and virulence-associated genomic and phenotypic
196 traits.

197

198 **The genomes of the STEC O181:H4 strain 17-07187 and the O104:H4 outbreak strain**
199 **mainly differ in their O-antigen gene clusters and mobile genetic elements**

200 The chromosomes (without plasmids) of the O181:H4 strain 17-07187 and the O104:H4
201 outbreak strain FWSEC0009 were nearly identical (~94.5% nucleotide identity). The most
202 striking difference between them were their respective O-antigen gene clusters (OAGC)
203 (Fig.2A, 2B). Although these two clusters were both situated in the same location in the
204 chromosome, between *galF* and *hisI* (Fig.2B), their gene contents and organization were very
205 different. Furthermore, their respective GC contents, 36.8% for O181 and 37% for O104,
206 differed from the chromosome GC content (~50.7%), highlighting the likely role of lateral gene
207 transfer in driving OAGC exchange.

208 14 potential prophage regions are present in the O181:H4 strain and 16 in the O104:H4
209 outbreak strain (Fig.2A, Table S3). 11 of the 14 prophages exhibited substantial sequence
210 identity (83-99.9%) with their O104:H4 counterparts and importantly, the *stx2*-encoding
211 prophages were nearly identical (~99% identity) (Fig.2C, Table S4, Fig.2C). Both *stx* phages
212 are inserted into the tryptophan repressor binding protein gene *wrbA*.

213

214 The genome of the O181:H4 strain included three plasmids of ~81kb, ~76kb, and ~63kb (Table
215 S2). The largest O181:H4 plasmid (plasmid 1) was an incompatibility group I1 (Incl1) plasmid
216 that showed only partial homology to the ESBL resistance plasmid of the O104:H4 outbreak
217 strain (~57% nucleotide identity) but was very similar (> 90% nucleotide identity) to
218 pHUSEC41-1 of STEC O104:H4 HUSEC41 from 2001 (92kb) (Fig.S1A) [3, 31]. Both of these
219 plasmids encode the *pill-V* genes for thin pili. Unlike pHUSEC041-1, the O181:H4 plasmid 1
220 did not contain antibiotic resistance genes (Fig.S1A, Tab.S5). O181:H4 plasmid 2 (76 kb) was
221 nearly identical (95% identity, 100% coverage) to the O104:H4 outbreak strain pAA (pAA-

222 EA11) and harbored virulence associated loci including *aggA/B/C/D*, which encode the AAF/I
223 fimbriae that promote bacterial adherence to host cells (Fig.1A, Fig.S1B, Tab.S6) [7, 10, 32].
224 O181:H4 plasmid 3 (63 kb) was not found in the O104:H4 outbreak strain; instead, it showed
225 similarity to DHA plasmids of several enterobacteria coding for AmpC β -lactamase [33].
226 However, unlike the DHA plasmids, resistance determinants were not present in the O181:H4
227 strain (Fig.S1C, Tab.S7). Together these observations reveal the striking similarity of the
228 chromosomes of the O181:H4 strain and the O104:H4 outbreak strain and that their chief
229 differences are confined to hot spots of recombination, i.e. their O-antigen gene clusters, and
230 to mobile genetic elements, particularly their plasmids.

231

232 **Additional recent global isolates of serotype O181:H4 and five other O groups belong** 233 **to ST678**

234 Next, we identified 158 genomes belonging to ST678 in Enterobase [24], which contains
235 ~202,200 *E. coli* genomes (as of April 11th, 2022). For a subset, serotype identification was
236 not available, and for these cases, serotype was predicted based on the available genomic
237 information[16]. 123 of the ST678 strains were O104:H4, however, 18 additional O181:H4
238 strains of ST678 were found. Furthermore, seven O127:H4, three O131:H4, and one each of
239 O69:H4 and OX13:H4 were identified (Fig.3A, Tab.S1). We categorized these as non-O104:H4
240 ST678 strains. Additionally, based on a close phylogenetic relationship and single difference
241 in MLST alleles, we added three non-ST678 strains to the non-O104:H4 ST678 category: two
242 O181:H4 strains (1472912 and 1472968) and one strain of the new and provisionally assigned
243 genoserotype OgN-RKI9:H4 (strain 608450) (Fig.3A, Tab.S1).

244 cgMLST confirmed the close phylogenetic relationship (allelic distance ~50) between AAF/I
245 gene positive O104:H4 strains including the outbreak strain (FWSEC0009), and the 21
246 O181:H4 strains, the OX13:H4 strain, and three of the seven O127:H4 strains (Fig.3A, 3B). All
247 of these strains showed AAF/I which were also found in the 2011 outbreak strain [5, 7]. Nine
248 strains (four O127:H4, three O131:H4, one each of O69:H4 and OgN-RKI9:H4) had higher
249 allelic distances up to 189. In contrast, these nine strains encoded AAF/III genes (Fig.3A, 3B).

250 The 25 AAF/I positive non-O104:H4 ST678 strains were all isolated in or after 2011 and were
251 associated with diarrheal disease and a subset of six O181:H4 strains harbored *stx2a* (Fig.3A,
252 3B, Fig.S2A, and Tab.S1). Interestingly, four of these six strains shared a very similar *stx*
253 phage with the O104:H4 outbreak strain (including strain 17-07187 from Germany), but two of
254 the six prophages were more distinct (Fig.3A). AAF/I-positive strains shared regions with a
255 higher relatedness than those positive for AAF/III with the pAA plasmid of the 2011 outbreak
256 strain.

257 The virulence gene profiles of the 34 non-O104 ST678 strains were generally similar to the
258 O104:H4 outbreak strain; however, there were a few differences. Specifically, *sepA* was
259 exclusively found in AAF/I-positive strains and EAST1 was present in all AAF/III-positive
260 strains and in only three of the AAF/I-positive isolates (Tab.S8).

261 The 34 non-O104:H4 ST678 strains were isolated in countries of Europe, Africa, and North
262 and South America. In addition, several were from individuals with a travel history that might
263 link these to East Asia (Fig.3A, S2B and C). Together, these observations show that ST678 *E.*
264 *coli* strains are found among seven different *E. coli* serotypes that have been linked to diarrheal
265 disease on several continents primarily during the past decade.

266

267 **Phylogenomic analyses of ST678 *E. coli* suggest recency of O-antigen gene exchange**

268 The O-antigen gene clusters encoding the six O groups in the 34 non-O104:H4 ST678 strains
269 are found in the same chromosomal location as the 2011 O104:H4 outbreak strain but are
270 composed of largely disparate genes (Fig.4A). These clusters also have a distinct GC content
271 (36.8-42.1%) from the backbone genome (50.5-50.7%) (Tab.S9), suggesting that they were
272 acquired by horizontal gene exchange. To explore the phylogenetic relationships among the
273 34 non-O104:H4 ST678 strains and a set of O104:H4 strains, their shared single nucleotide
274 polymorphisms (SNPs) were analyzed using the O104:H4 outbreak strain FWSEC0009 as a
275 reference. A positive correlation ($R=0.81$; $R^2=0.66$) was found between isolation time and
276 genetic divergence (Fig.S3) which supports that mutations have accumulated in a clock-like
277 fashion without notable outliers. Fig.4B shows that the AAF/III-positive strains (upper part),

278 which include nine non-O104:H4 strains of four different serotypes, are clearly separated from
279 the AAF/I positive strains (lower part, referred to as clade I). Clade I is comprised of O104:H4
280 strains, such as from the 2011 outbreak, 21 O181:H4 strains, three O127:H4 strains, and an
281 OX13:H4 strain (Fig.4B). Strains of clade I are much more closely related to one another than
282 to the AAF/III positive strains. The structure of the phylogeny suggests that clade I strains are
283 derived from an AAF/III-positive precursor. Within clade I, two subclades, Ia and Ib, contain
284 non-O104:H4 strains. Clade Ia is composed of O104:H4 non-outbreak strains isolated from
285 2015-2021 in the United Kingdom and in Kenya and the 2018 OX13:H4 strain that was
286 associated with travel to Ethiopia. Clade Ib contains the 21 O181:H4 strains and three O127:H4
287 strains isolated from 2011 to 2021 from diverse continents, but in contrast to the other branches
288 in clade I, no O104:H4 strains belong to this subclade. In the phylogeny, the branch containing
289 clade Ib arises from an older AAF/I positive O104:H4 dominated branch, suggesting these
290 serotypes are derived from an O104:H4 precursor (Fig.4B). Similarly, the phylogeny suggests
291 that the AAF/I positive O127:H4 strains of clade Ib arose from an O181:H4 precursor. Together
292 these observations suggest that O-antigen gene cluster exchange has occurred repeatedly
293 within ST678 strains. This is illustrated by the appearance of O127:H4 strains at multiple
294 locations in the tree and suggests that O127 O-antigen gene donor strains share a niche with
295 ST678 recipient strains.

296

297 **Discussion**

298 The *E. coli* O104:H4 outbreak in Germany in the early summer of 2011 was a public health
299 emergency; however, this serotype was rarely isolated as a cause of HUS after the epidemic
300 subsided. Nevertheless, our findings suggest that the unusual set of virulence factors that
301 characterized this Shiga toxin 2-producing EAEC strain remains a threat to human health.
302 Serotype conversion has cloaked this highly virulent genotype with several O-antigens,
303 including O181, O127, and OX13, which are present both in *stx* positive and *stx* negative
304 disease-linked isolates closely related to the O104:H4 outbreak strain. We found non-O104:H4
305 ST678 strains from a variety of countries in Europe, Africa, and the Americas and several of

306 the cases were associated with travel to Africa and Asia, suggesting these virulence-
307 associated strains are globally distributed. Like the O104:H4 outbreak strain, the AAF/I positive
308 ST678 strains had similar chromosomes, similar pAA-linked virulence genes, as well as
309 virulence factors including the SPATEs, SigA, and Pic [5, 7, 8, 11]. The most salient difference
310 of the chromosomes of these strains with that of the outbreak strain were their respective O-
311 antigen gene clusters. Thus, horizontal exchange of these clusters appears to have been a
312 critical step in the evolution of these new pathogenic serotypes, some of which were linked to
313 HUS or bloody diarrhea.

314

315 The prime example uncovered here is the derivation of O181:H4 pathogens from an O104:H4
316 precursor via O-antigen gene exchange. Among the 21 O181:H4 strains, six harbored a *stx2a*-
317 encoding prophage, including HUS-linked strain 17-07187. In four of these strains, the *stx2a*
318 prophage was very similar to the *stx2a* prophage of the 2011 outbreak strain, suggesting that
319 the O-antigen exchange was a more recent event in their evolution compared to the acquisition
320 of the Shiga toxin-encoding prophage (Fig 3A). The absence of *stx* prophages in 15 of the
321 O181:H4 isolates may be due to a lack of *stx* prophage acquisition or have resulted from loss
322 of their *stx* prophages, which is well documented in STEC in other serotypes [34]. Thus, in
323 addition to O-antigen exchange, on-going horizontal transmission of mobile genetic elements,
324 such as *stx* phages, has contributed to diversification of diarrheagenic ST678.

325

326 O-antigen gene clusters of Gram-negative bacteria are hot spots for diversifying selection and
327 recombination events [15, 19]. Serotype conversion by lateral exchange of OAGCs has played
328 an important role in the evolution of enteric pathogens. For example, *Vibrio cholerae* serogroup
329 O139, which arose by exchange of O1 and O139 OAGCs, transiently replaced the dominant
330 O1 group as the cause of cholera in 1992-93 [13, 35]. Also, an interspecies exchange was
331 described for O8 and O9 O-antigens of *E. coli*, which are identical to O5 and O3 of *Klebsiella*
332 *pneumoniae*, respectively [36]. Our discovery of seven distinct O-groups that share the
333 flagellar H4 antigen and a highly similar virulence-linked genetic makeup provides a compelling

334 example of the role of O-antigen exchange in the diversification of diarrheagenic *E. coli*. In
335 addition, STEC O157:H7 is thought to have arisen from an enteropathogenic *E. coli* (EPEC)-
336 like O55:H7 strain which initially acquired *stx* via phage transduction and subsequently
337 acquired the O157 O group by exchange of the O55 with the O157 OAGC [37]. Also, an O182-
338 O156 switch is thought to have occurred in STEC strains persistently infecting cattle [18].
339 Consequently, different OAGs may be found in highly related genotypes [38, 39].

340

341 We can only speculate about the conditions driving OAGC exchange among ST678 strains.
342 Shiga toxin producing O104:H4 EAEC strains, such as the 2011 outbreak strain, are
343 considered human-restricted pathogens and have not been isolated from animals, such as
344 cattle [40]. It is possible that the OAGC exchange occurred in a human host, where the human
345 intestinal microbiome may contain *E. coli* strains of O groups, such as O181 and O127, that
346 could have donated their OAGC to an O104:H4 recipient by some means of horizontal
347 exchange. Also, epidemiological investigations suggest that fenugreek sprouts were the food
348 source that initiated the STEC O104:H4 epidemic; therefore plants colonized by microbiota
349 may be another possible site for OAGC exchange [4].

350

351 In conclusion, our study highlights how lateral exchange of O-antigen gene clusters can lead
352 to the rapid diversification of a globally important pathogen. Surveillance to uncover how highly
353 virulent strains may reemerge and spread in new O-antigen outfits is warranted. Further, an
354 important clinical implication of these findings is that serotype identification cannot be used as
355 a simple proxy for strain virulence and needs to be complemented by comprehensive virulence
356 gene analysis.

357

358 **Footnotes:**

359 **Funding**

360 The study sponsors had no role in study design; collection, analysis, and interpretation of data;
361 writing of the manuscript; and in the decision to submit it for publication.

362

363 **Acknowledgements**

364 We thank Ute Strutz, Thomas Garn, Tobias Größl, and Karsten Großhennig for excellent
365 technical assistance (Robert Koch Institute, Wernigerode, Germany). We also acknowledge
366 the MF1 bioinformatics and MF2 genome sequencing unit of the Robert Koch Institute (Berlin)
367 for the support of bioinformatic analysis and Illumina MiSeq and MinION sequencing. We thank
368 members of the Waldor and Flieger groups for helpful comments and discussions during the
369 study. AFI received funding from the German Ministry of Health (IGS-ZOO grant number
370 2518FSB706) and AFI and AFr from the intensified molecular surveillance initiative of the
371 Robert Koch- Institute. IWC received funding from the National Institute of Diabetes and
372 Digestive and Kidney Diseases of the National Institutes of Health (Award Number 5
373 T32DK007477-37).

374

375 **Conflict of interests**

376 We declare no conflict of interests.

377

378 **Data availability**

379 All data generated or analyzed during this study are available in this Article and the
380 appendices. The sequences were uploaded to NCBI project: PRJNA833419.

381

382 **Corresponding author contact information:**

383 Antje Flieger

384 email: fliegera@rki.de

385 Tel.: +49-30-18754-2522

386 **References**

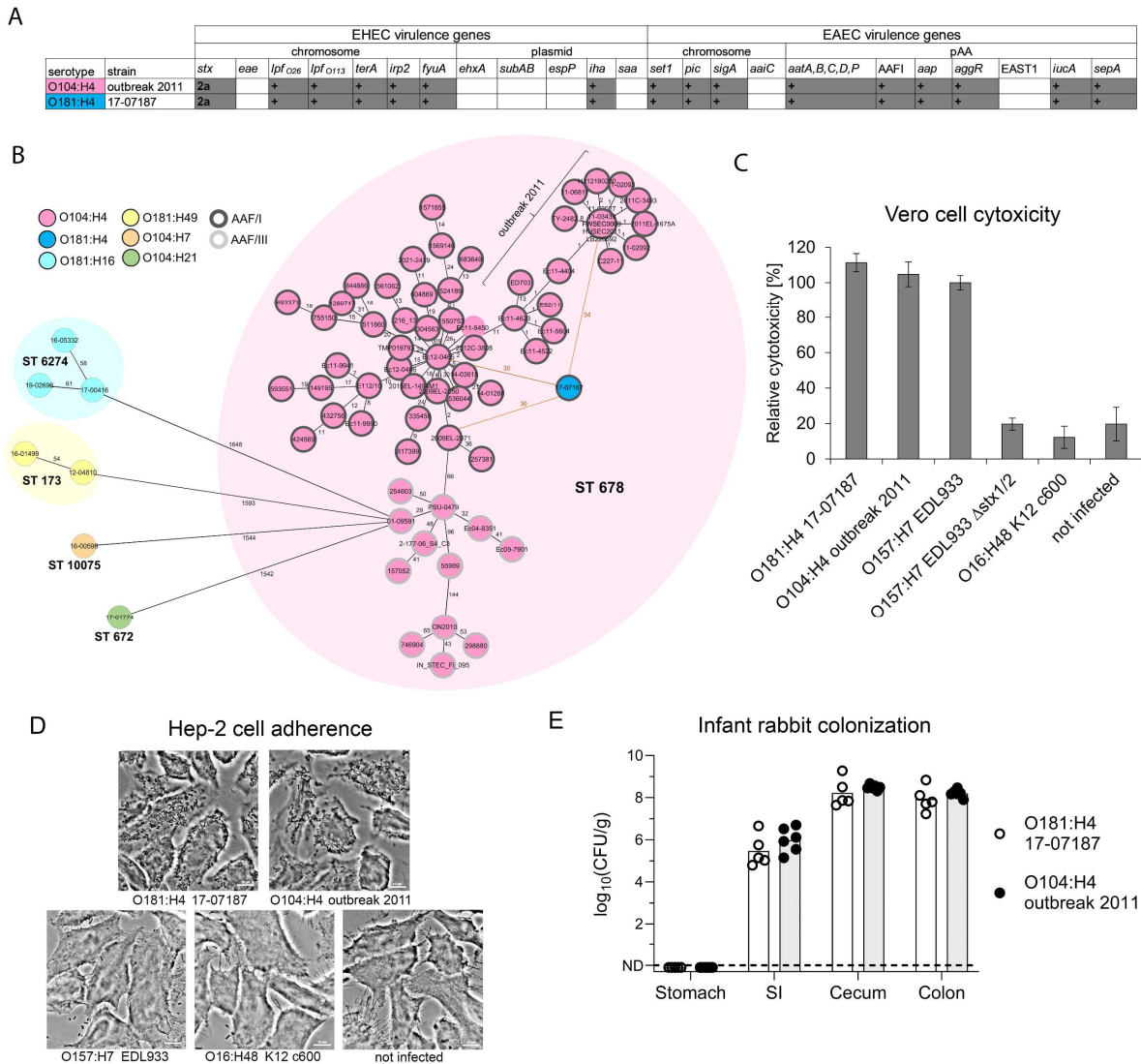
- 387 1. Croxen MA, Law RJ, Scholz R, Keeney KM, Wlodarska M, Finlay BB. Recent advances in
388 understanding enteric pathogenic *Escherichia coli*. Clin Microbiol Rev **2013**; 26:822-80.
- 389 2. Caprioli A, Scavia G, Morabito S. Public Health Microbiology of Shiga Toxin-Producing
390 *Escherichia coli*. Microbiol Spectr **2014**; 2:EHEC-0014-2013.
- 391 3. Mellmann A, Bielaszewska M, Kock R, et al. Analysis of collection of hemolytic uremic
392 syndrome-associated enterohemorrhagic *Escherichia coli*. Emerg Infect Dis **2008**; 14:1287-
393 90.
- 394 4. Werber D, Krause G, Frank C, et al. Outbreaks of virulent diarrheagenic *Escherichia coli*--
395 are we in control? BMC Med **2012**; 10:11.
- 396 5. Bielaszewska M, Mellmann A, Zhang W, et al. Characterisation of the *Escherichia coli*
397 strain associated with an outbreak of haemolytic uraemic syndrome in Germany, 2011: a
398 microbiological study. Lancet Infect Dis **2011**; 11:671-6.
- 399 6. Frank C, Werber D, Cramer JP, et al. Epidemic profile of Shiga-toxin-producing
400 *Escherichia coli* O104:H4 outbreak in Germany. The New England journal of medicine **2011**;
401 365:1771-80.
- 402 7. Rasko DA, Webster DR, Sahl JW, et al. Origins of the *E. coli* strain causing an outbreak of
403 hemolytic-uremic syndrome in Germany. The New England journal of medicine **2011**;
404 365:709-17.
- 405 8. Rohde H, Qin J, Cui Y, et al. Open-source genomic analysis of Shiga-toxin-producing *E.*
406 *coli* O104:H4. The New England journal of medicine **2011**; 365:718-24.
- 407 9. Mellmann A, Harmsen D, Cummings CA, et al. Prospective genomic characterization of
408 the German enterohemorrhagic *Escherichia coli* O104:H4 outbreak by rapid next generation
409 sequencing technology. PLoS One **2011**; 6:e22751.
- 410 10. Boisen N, Melton-Celsa AR, Scheutz F, O'Brien AD, Nataro JP. Shiga toxin 2a and
411 Enteroaggregative *Escherichia coli*--a deadly combination. Gut Microbes **2015**; 6:272-8.

- 412 11. Munera D, Ritchie JM, Hatzios SK, et al. Autotransporters but not pAA are critical for
413 rabbit colonization by Shiga toxin-producing *Escherichia coli* O104:H4. *Nat Commun* **2014**;
414 5:3080.
- 415 12. Coipan CE, Friesema IH, van den Beld MJC, et al. Sporadic Occurrence of
416 Enteroaggregative Shiga Toxin-Producing *Escherichia coli* O104:H4 Similar to 2011
417 Outbreak Strain. *Emerg Infect Dis* **2022**; 28:1890-4.
- 418 13. Lerouge I, Vanderleyden J. O-antigen structural variation: mechanisms and possible
419 roles in animal/plant-microbe interactions. *FEMS Microbiol Rev* **2002**; 26:17-47.
- 420 14. Orskov F, Orskov I. *Escherichia coli* serotyping and disease in man and animals. *Can J*
421 *Microbiol* **1992**; 38:699-704.
- 422 15. Iguchi A, Iyoda S, Kikuchi T, et al. A complete view of the genetic diversity of the
423 *Escherichia coli* O-antigen biosynthesis gene cluster. *DNA Res* **2015**; 22:101-7.
- 424 16. Lang C, Hiller M, Konrad R, Fruth A, Flieger A. Whole-Genome-Based Public Health
425 Surveillance of Less Common Shiga Toxin-Producing *Escherichia coli* Serovars and
426 Untypeable Strains Identifies Four Novel O Genotypes. *J Clin Microbiol* **2019**; 57.
- 427 17. Samuel G, Hogbin JP, Wang L, Reeves PR. Relationships of the *Escherichia coli* O157,
428 O111, and O55 O-antigen gene clusters with those of *Salmonella enterica* and *Citrobacter*
429 *freundii*, which express identical O antigens. *J Bacteriol* **2004**; 186:6536-43.
- 430 18. Geue L, Menge C, Eichhorn I, et al. Evidence for Contemporary Switching of the O-
431 Antigen Gene Cluster between Shiga Toxin-Producing *Escherichia coli* Strains Colonizing
432 Cattle. *Front Microbiol* **2017**; 8:424.
- 433 19. Milkman R, Jaeger E, McBride RD. Molecular evolution of the *Escherichia coli*
434 chromosome. VI. Two regions of high effective recombination. *Genetics* **2003**; 163:475-83.
- 435 20. Prager R, Lang C, Aurass P, Fruth A, Tietze E, Flieger A. Two novel EHEC/EAEC hybrid
436 strains isolated from human infections. *PLoS One* **2014**; 9:e95379.
- 437 21. Walker BJ, Abeel T, Shea T, et al. Pilon: an integrated tool for comprehensive microbial
438 variant detection and genome assembly improvement. *PLoS One* **2014**; 9:e112963.

- 439 22. Darling AE, Mau B, Perna NT. progressiveMauve: multiple genome alignment with gene
440 gain, loss and rearrangement. PLoS One **2010**; 5:e11147.
- 441 23. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
442 improvements in performance and usability. Mol Biol Evol **2013**; 30:772-80.
- 443 24. Zhou Z, Alikhan NF, Mohamed K, Fan Y, Agama Study G, Achtman M. The EnteroBase
444 user's guide, with case studies on Salmonella transmissions, *Yersinia pestis* phylogeny, and
445 *Escherichia* core genomic diversity. Genome Res **2020**; 30:138-52.
- 446 25. Arndt D, Grant JR, Marcu A, et al. PHASTER: a better, faster version of the PHAST
447 phage search tool. Nucleic Acids Res **2016**; 44:W16-21.
- 448 26. Aziz RK, Bartels D, Best AA, et al. The RAST Server: rapid annotations using
449 subsystems technology. BMC Genomics **2008**; 9:75.
- 450 27. Lang C, Fruth A, Holland G, et al. Novel type of pilus associated with a Shiga-toxigenic *E.*
451 *coli* hybrid pathovar conveys aggregative adherence and bacterial virulence. Emerg
452 Microbes Infect **2018**; 7:203.
- 453 28. Croucher NJ, Page AJ, Connor TR, et al. Rapid phylogenetic analysis of large samples of
454 recombinant bacterial whole genome sequences using Gubbins. Nucleic Acids Res **2015**;
455 43:e15.
- 456 29. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms
457 and methods to estimate maximum-likelihood phylogenies: assessing the performance of
458 PhyML 3.0. Syst Biol **2010**; 59:307-21.
- 459 30. Rambaut A, Lam TT, Max Carvalho L, Pybus OG. Exploring the temporal structure of
460 heterochronous sequences using TempEst (formerly Path-O-Gen). Virus Evol **2016**;
461 2:vew007.
- 462 31. Kunne C, Billion A, Mshana SE, et al. Complete sequences of plasmids from the
463 hemolytic-uremic syndrome-associated *Escherichia coli* strain HUSEC41. J Bacteriol **2012**;
464 194:532-3.

- 465 32. Schiller P, Knodler M, Berger P, et al. The Superior Adherence Phenotype of *E. coli*
466 O104:H4 is Directly Mediated by the Aggregative Adherence Fimbriae Type I. *Virulence*
467 **2021**; 12:346-59.
- 468 33. Gaillot O, Clement C, Simonet M, Philippon A. Novel transferable beta-lactam resistance
469 with cephalosporinase characteristics in *Salmonella enteritidis*. *J Antimicrob Chemother*
470 **1997**; 39:85-7.
- 471 34. Bielaszewska M, Prager R, Kock R, et al. Shiga toxin gene loss and transfer in vitro and
472 in vivo during enterohemorrhagic *Escherichia coli* O26 infection in humans. *Applied and*
473 *environmental microbiology* **2007**; 73:3144-50.
- 474 35. Faruque SM, Sack DA, Sack RB, Colwell RR, Takeda Y, Nair GB. Emergence and
475 evolution of *Vibrio cholerae* O139. *Proc Natl Acad Sci U S A* **2003**; 100:1304-9.
- 476 36. Sugiyama T, Kido N, Kato Y, Koide N, Yoshida T, Yokochi T. Evolutionary relationship
477 among rfb gene clusters synthesizing mannose homopolymer as O-specific polysaccharides
478 in *Escherichia coli* and *Klebsiella*. *Gene* **1997**; 198:111-3.
- 479 37. Feng P, Lempel KA, Karch H, Whittam TS. Genotypic and phenotypic changes in the
480 emergence of *Escherichia coli* O157:H7. *J Infect Dis* **1998**; 177:1750-3.
- 481 38. Ingle DJ, Valcanis M, Kuzevski A, et al. In silico serotyping of *E. coli* from short read data
482 identifies limited novel O-loci but extensive diversity of O:H serotype combinations within and
483 between pathogenic lineages. *Microb Genom* **2016**; 2:e000064.
- 484 39. Ochman H, Selander RK. Evidence for clonal population structure in *Escherichia coli*.
485 *Proc Natl Acad Sci U S A* **1984**; 81:198-201.
- 486 40. Wieler LH, Semmler T, Eichhorn I, et al. No evidence of the Shiga toxin-producing *E. coli*
487 O104:H4 outbreak strain or enteroaggregative *E. coli* (EAEC) found in cattle faeces in
488 northern Germany, the hotspot of the 2011 HUS outbreak area. *Gut Pathog* **2011**; 3:17.
- 489
- 490

491 **Figures**



492

493 **Figure 1: STEC O181:H4/17-07187 shares close phylogenetic relationship and similar**

494 **virulence traits with O104:H4 outbreak strain. (A) Overlapping virulence gene profiles of**

495 **the O104:H4 outbreak strain (FWSEC0009) and the O181:H4 strain 17-07187. (B) Minimal**

496 **spanning tree of O181:H4 17-07187 and selected O104:H4 strains (representing phylogenetic**

497 **diversity of O104:H4 strains with respect to isolation time from 1998 to 2022 and location)**

498 **based on cgMLST involving 2513 alleles confirms their close phylogenetic relationship.**

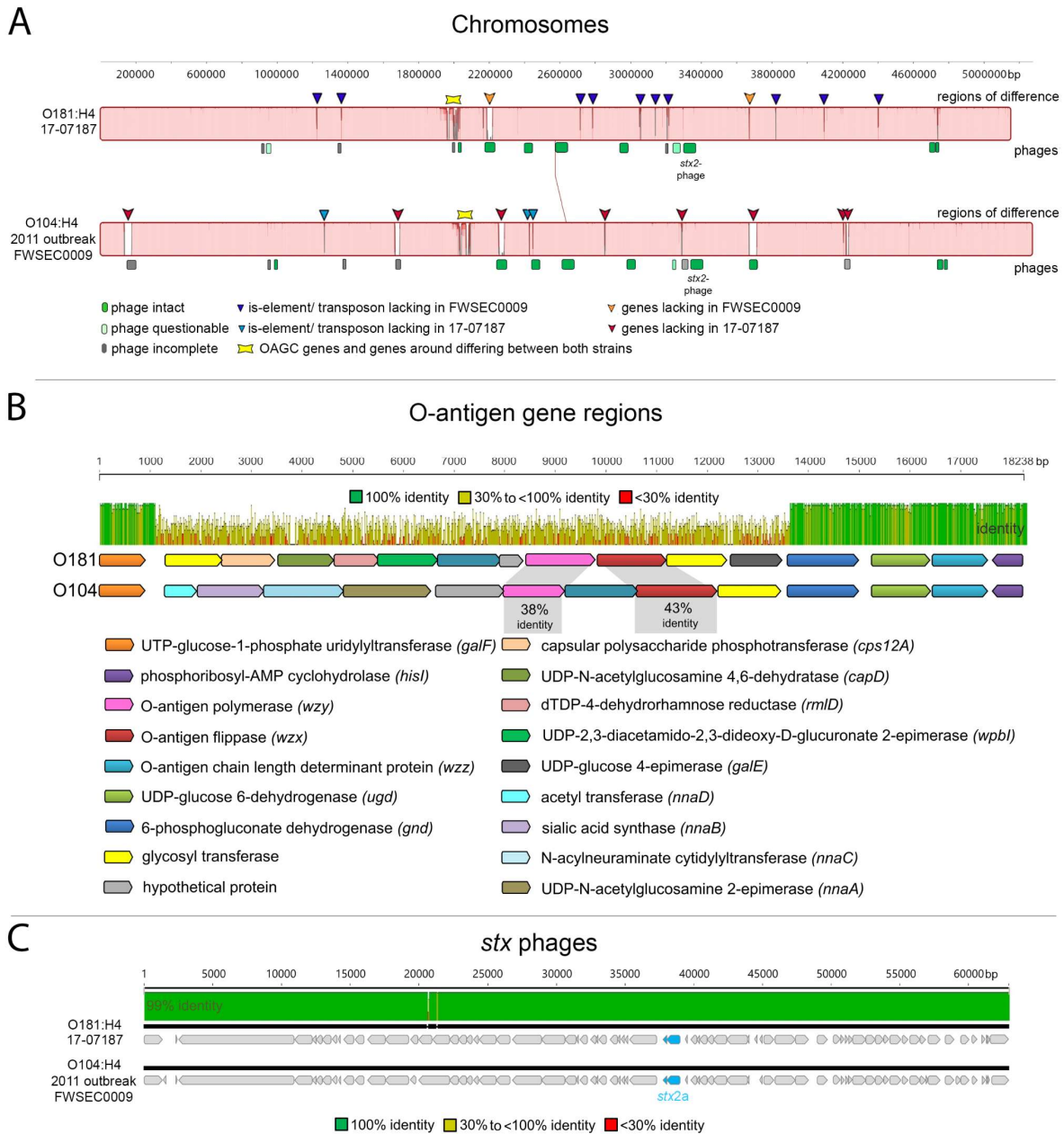
499 **O181:H4 17-07187 and O104:H4 strains 2009EL-2071, Ec12-0465, and outbreak strains from**

500 **2011 (including strains FWSEC0009, C227-11, and 11-02027) differ by only 30-34 alleles.**

501 **Serotypes and AAF type are indicated in the key. (C) Cytotoxicity of O181:H4 17-07187**

502 towards Vero cells and (D) adherence pattern to HEp-2 cells are comparable to the O104:H4
503 outbreak strain 11-02027. Results are representative of at least two additional experiments.
504 EHEC O157:H7 EDL933 producing Stx1 and Stx2 served as reference in the cytotoxicity assay
505 and was set to 100%. The cytotoxicity results represent the means and standard deviations of
506 triplicate samples (n=3) and are representative of at least two additional experiments. E)
507 Intestinal colonization in infant rabbits inoculated with O104:H4 outbreak strain C227-11 (n=6)
508 or O181:H4 strain 17-07187 (n=5). CFU/g (colony forming units per gram of tissue) denotes
509 concentration of bacteria recovered three days post inoculation from homogenated intestinal
510 tissues. Data points represent individual rabbits from two independent litters split between the
511 two strains. Bars show the geometric mean, SI is small intestine, ND is not detected.

512



513

514 **Figure 2: O181:H4 17-07187 and O104:H4 outbreak strain genomes mainly differ in**

515 **OAGCs and mobile genetic elements.** (A) Whole chromosome MAUVE alignment of

516 O181:H4 17-07187 and the O104:H4 outbreak strain FWSEC0009 highlighting mobile genetic

517 elements and differences in phage regions, IS-/ transposon elements, and OAGC regions. (B)

518 OAGCs of O181:H4 17-07187 and the O104:H4 FWSEC0009 outbreak strain are flanked by

519 homologous upstream (*galF*) and downstream (*gnd* to *hisI*) regions. MAFFT alignment shows

520 that the regions in between *galF* and *gnd* are very different. (C) MAFFT alignment of the *stx2a*-

521 encoding phages in O181:H4 17-07187 and the outbreak strain O104:H4 outbreak strain

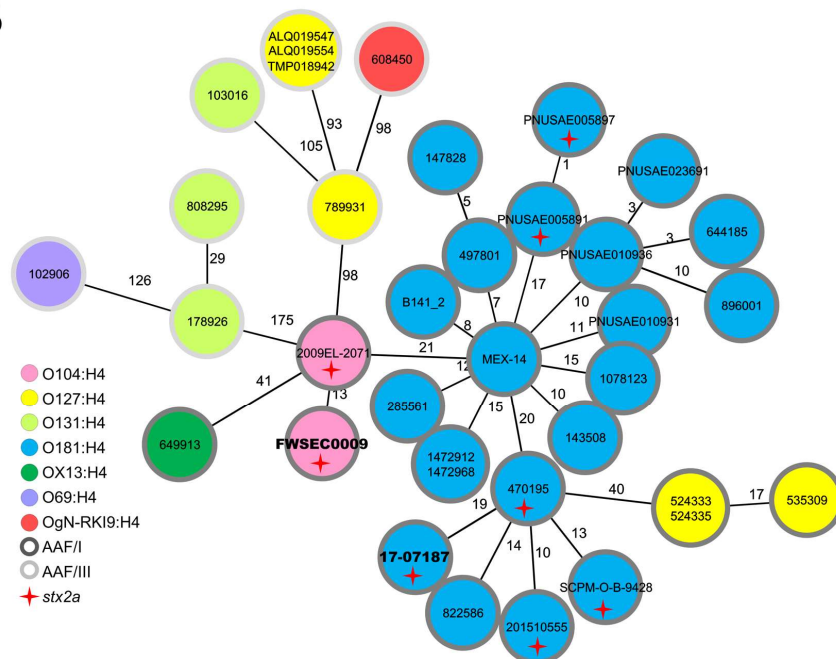
522 FWSEC0009 shows that they are very similar (99% nucleotide identity).

A

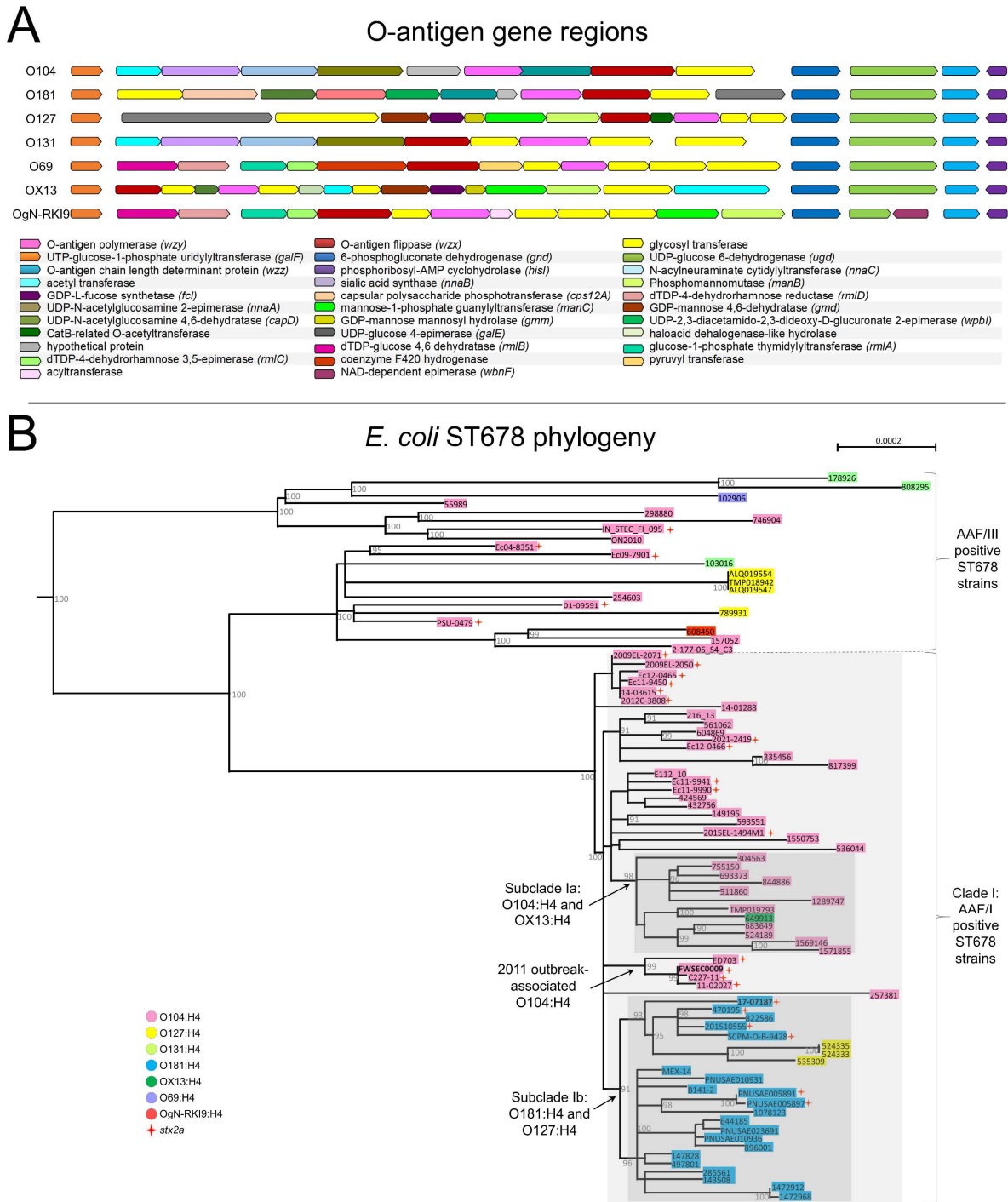
strain	serotype	allele distance to strain FWSEC0009	country of isolation	travel association	year of isolation	FWSEC0009 O104:H4			
						p1	p2 pAA	p3	stx2a-phage
FWSEC0009	O104:H4	0	Germany		2011	100	100	100	100
MEX-14	O181:H4	26	Mexico		2011	5.7	100	100	
497801	O181:H4	27	United Kingdom	Morocco	2013	1.7	100	100	
147828	O181:H4	28	United Kingdom	Morocco	2015	90.7	100	100	
B141_2	O181:H4	28	Ecuador		2015	2.9	95.5	100	
143508	O181:H4	30	United Kingdom		2015	17.5	100	100	
PNUSAE010936	O181:H4	30	United States		2017	1.1	100	100	
644185	O181:H4	31	United Kingdom	Mexico	2018	8.4	100	100	
PNUSAE023691	O181:H4	31	United States		2019	1.1	100	99.9	
285561	O181:H4	32	United Kingdom	Egypt	2016	32.1	100	100	
470195	O181:H4	32	United Kingdom	India	2018	16.1	100	20.3	100
PNUSAE010931	O181:H4	32	United States		2017	1.6	100	100	
201510555	O181:H4	34	France		2015	9	100	0	100
1078123	O181:H4	35	United Kingdom		2021	1.9	97.1	100	
1472912	O181:H4	35	United Kingdom		2021	4.7	100	100	
17-07187	O181:H4	35	Germany		2017	37.3	100	0	100
1472968	O181:H4	36	United Kingdom		2021	4.3	100	100	
SCPM-O-B-9428	O181:H4	36	Russia		2018	36.8	100	87.8	98.9
822586	O181:H4	38	United Kingdom	Morocco	2019	3.3	99.7	22.1	
896001	O181:H4	38	United Kingdom		2020	32.7	99.6	100	
PNUSAE005891	O181:H4	38	United States		2016	35.6	100	100	58.1
PNUSAE005897	O181:H4	39	United States		2017	35	99.9	100	58.4
649913	OX13:H4	46	United Kingdom	Ethiopia	2018	5.4	99.4	100	
524333	O127:H4	53	United Kingdom	India	2018	3.6	100	18.9	
524335	O127:H4	53	United Kingdom	India	2018	3	100	53.6	
535309	O127:H4	54	United Kingdom		2018	14.8	100	19.4	
789931	O127:H4	106	United Kingdom	Kenya	2019	17.5	70.8	13	
608450	RK19:H4	112	United Kingdom		2018	7.3	62.1	30.5	
ALQ019547	O127:H4	113	Kenya		2015	16.2	69.6	0	
ALQ019554	O127:H4	113	Kenya		2015	16.4	69.3	0	
TMP018942	O127:H4	113	Kenya		2015	16.1	69.5	0	
103016	O131:H4	121	Gambia		2015	3.5	60.4	5.9	
178926	O131:H4	182	United Kingdom	India	2019	31.2	62.9	13	
102906	O69:H4	187	Gambia		2009	15	70.2	5.7	
808295	O131:H4	189	United Kingdom	India	2019	4.8	60.9	52.6	



B



524 **Figure 3: Recent global isolates of serotype O181:H4 and five other O groups belong to**
525 **ST678.** (A) Summary of serotypes, isolation dates, along with relatedness of genomes
526 (cgMLST-based allelic distances) and mobile genetic elements (plasmids and *stx* phage) of
527 the 34 non-O104:H4 ST678 (ST12598, ST12610) clinical strains (found on Enterobase)
528 compared to the O104:H4 outbreak strain FSWEC0009. (B) Serotypes of the non-O104:H4
529 ST678 strains and the O104:H4 strains 2009EL-207 and the outbreak strain FSWEC0009 in a
530 minimal spanning tree based on cgMLST. Although, phylogenetically very close to ST678
531 strains, O181:H4 strains 1472912 and 1472968 in fact belong to ST12610 due to a point
532 mutation in *icd* and OgN-RKI9:H4 strain 608450 to ST12598 due to a point mutation in *recA*.
533



541 rooted with an outgroup consisting of *E. coli* strains K12 c600, EcO42, EDL933, O157:H7 strain
542 Sakai, O26:H11 strain 11368, and O103:H2 strain 12009. Bootstrap values >90% are
543 indicated. Serotypes are as indicated in the key and presence of *stx* is indicated by a red star.
544 Scale bar refers to a phylogenetic distance of 0.0002 nucleotide substitutions per site.
545

546 **Supplementary Data:**

547 **Figure S1:** The O181:H4 STEC strain 17-07187 and the O104:H4 outbreak strain harbor a
548 partially distinct set of plasmids with overlapping pAA. MAUVE alignments highlight that (A)
549 STEC O181:H4 17-07187 plasmid 1 is most similar to pHUSEC41-1 of STEC O104:H4
550 HUSEC41 from 2001 but lacks antibiotic resistance genes, (B) O181:H4 17-07187 plasmid 2
551 is very similar to pAA of the O104:H4 outbreak strain and (C) O181:H4 17-07187 plasmid 3 is
552 similar to DHA-1 plasmids of enterobacteria but lacks antibiotic resistance genes.

553

554 **Figure S2:** Minimal spanning tree based on cgMLST of the 34 non-O104:H4 ST678 strains
555 and the O104:H4 strains 2009EL-2071 and FWSEC0009 colored by year of isolation (A),
556 country of isolation (B), and travel association (C).

557

558 **Figure S3:** TempEst revealed a positive correlation ($R=0.81$; $R^2=0,66$) between isolation time
559 and genetic divergence represented by the root to tip regression analysis using RAxML tree
560 generated by Gubbins based on a recombination-corrected alignment of genome wide
561 polymorphic sites with O104:H4 2011 outbreak strain FWSEC0009 as reference.

562

563 **Table S1:** Overview of all strains used in the study and information on their serotypes,
564 virulence genes, country of isolation, travel association, clinical data, and genome
565 accessions. (provided as excel file)

566

567 **Table S2:** Genome characteristics of O181:H4 17-07187 and the 2011 O104:H4 outbreak
568 strain FWSEC0009. (provided as word file)

569

570 **Table S3:** Comparison of phage-regions detected in STEC O181:H4 17-07187 and in STEC
571 O104:H4 strains FWSEC00009 from the outbreak 2011 and 2009EL-2071 from 2009.
572 (provided as excel file)

573

574 **Table S4:** Comparison of the *stx2a* phages detected in STEC O181:H4 17-07187 and in the
575 O104:H4 outbreak strain FWSEC00009. (provided as excel file)

576

577 **Table S5:** Comparison of O181:H4 17-07187 plasmid 1 to pHUSEC41-1 of STEC O104:H4
578 HUSEC41 from 2001. (provided as excel file)

579

580 **Table S6:** Comparison of O181:H4 17-07187 plasmid 2 to pAA of the O104:H4 outbreak
581 strain. (provided as excel file)

582

583 **Table S7:** Comparison of O181:H4 17-07187 plasmid 3 to pDHA1 of enterobacteria.
584 (provided as excel file)

585 **Table S8:** Overview of the virulence gene profiles of the 34 non-O104:H4 ST678 strains and
586 the 2011 O104:H4 outbreak strain FWSEC00009. (provided as excel file)

587

588 **Table S9:** GC contents of *E. coli* chromosomes and their respective O-antigen gene clusters.
589 (provided as excel file)