

A peer-reviewed version of this preprint was published in PeerJ on 17 February 2015.

[View the peer-reviewed version](https://peerj.com/articles/739) (peerj.com/articles/739), which is the preferred citable publication unless you specifically need to cite this preprint.

Ashton PM, Perry N, Ellis R, Petrovska L, Wain J, Grant KA, Jenkins C, Dallman TJ. 2015. Insight into Shiga toxin genes encoded by *Escherichia coli* O157 from whole genome sequencing. PeerJ 3:e739
<https://doi.org/10.7717/peerj.739>

Insight into Shiga toxin genes encoded by *Escherichia coli* O157 from whole genome sequencing

Philip Ashton, Neil Perry, Richard J Ellis, Liljana Petrovska, John Wain, Kathie Grant, Claire Jenkins, Tim Dallman

The ability of Shiga toxin-producing *Escherichia coli* (STEC) to cause severe illness in humans is determined by multiple host factors and bacterial characteristics, including Shiga toxin (Stx) subtype. Given the link between Stx2a subtype and disease severity, we sought to identify the *stx* subtypes present in whole genome sequences (WGS) of 444 isolates of STEC O157. Difficulties in assembling the *stx* genes in some strains, were overcome by using two complementary bioinformatics methods; mapping and *de novo* assembly. We compared the WGS analysis with the results obtained using a PCR approach and investigated the diversity within and between the subtypes. All strains of STEC O157 in this study had *stx1a*, *stx2a* or *stx2c* or a combination of these three genes. There was over 99% (442/444) concordance between PCR and WGS. When common source strains were excluded, 236/349 strains of STEC O157 had multiple copies of different Stx subtypes and 54 had multiple copies of the same Stx subtype. Of those strains harbouring multiple copies of the same Stx subtype, 33 had variants between the alleles while 21 had identical copies. Strains harbouring Stx2a only were most commonly found to have multiple alleles of the same subtype (42%). Both the PCR and WGS approach to *stx* subtyping provided a good level of sensitivity and specificity. In addition, the WGS data also showed there were a significant proportion of strains harbouring multiple alleles of the same Stx subtype associated with clinical disease in England.

2 **Insight into Shiga toxin genes encoded by *Escherichia coli* O157**
3 **from whole genome sequencing**

4 **Authors:**

5 Philip M. Ashton¹, Neil Perry¹, Richard Ellis², Liljana Petrovska², John Wain³, Kathie A. Grant¹,
6 Claire Jenkins¹, Tim J. Dallman¹#

7 **Affiliations:**

8 ¹Gastrointestinal Bacteria Reference Unit, Public Health England, 61 Colindale Avenue, London,
9 NW9 5HT

10 ²Animal Health and Veterinary Laboratories Agency, Woodham Lane, New Haw, Addlestone,
11 Surrey, KT15 3NB

12 ³University of East Anglia, Norwich Research Park, Norwich, NR4 7TJ

13 **Corresponding author:**

14 Dr Tim J. Dallman, Gastrointestinal Bacteria Reference Unit,
15 Public Health England, 61 Colindale Avenue, London, NW9 5EQ

16 Email: tim.dallman@phe.gov.uk

17 Tel : 0208 327 6420

18 Fax : 0208 327 7112

19 **Author emails**

20 Philip Ashton – philip.ashton@phe.gov.uk

21 Neil Perry – neil.perry@phe.gov.uk

22 Richard Ellis - Richard.Ellis@ahvla.gsi.gov.uk

23 Liljana Petrovska - Liljana.Petrovska@ahvla.gsi.gov.uk

24 John Wain - J.Wain@uea.ac.uk

25 Kathie Grant - Kathie.Grant@phe.gov.uk

26 Claire Jenkins - Claire.Jenkins1@phe.gov.uk

27 Tim Dallman - tim.dallman@phe.gov.uk

28 **Author contributions**

29 Conceived and designed the experiments: PA, TD, CJ

30 Performed the experiments: PA, NP, RE

31 Analysed the data: PA

32 Wrote the paper: PA, TD, CJ

33 Contributed reagents/materials/analysis tools: LP, JW, KG

34 **Running title:** Use of WGS for subtyping stx genes

35 **Key words:** STEC O157; WGS; stx subtyping

36 **Abstract**

37 The ability of Shiga toxin-producing *Escherichia coli* (STEC) to cause severe illness in humans is
38 determined by multiple host factors and bacterial characteristics, including Shiga toxin (Stx)
39 subtype. Given the link between Stx2a subtype and disease severity, we sought to identify the
40 *stx* subtypes present in whole genome sequences (WGS) of 444 isolates of STEC O157.
41 Difficulties in assembling the *stx* genes in some strains, were overcome by using two
42 complementary bioinformatics methods; mapping and *de novo* assembly. We compared the
43 WGS analysis with the results obtained using a PCR approach and investigated the diversity
44 within and between the subtypes. All strains of STEC O157 in this study had *stx1a*, *stx2a* or
45 *stx2c* or a combination of these three genes. There was over 99% (442/444) concordance
46 between PCR and WGS. When common source strains were excluded, 236/349 strains of STEC
47 O157 had multiple copies of different Stx subtypes and 54 had multiple copies of the same Stx
48 subtype. Of those strains harbouring multiple copies of the same Stx subtype, 33 had variants
49 between the alleles while 21 had identical copies. Strains harbouring Stx2a only were most
50 commonly found to have multiple alleles of the same subtype (42%). Both the PCR and WGS
51 approach to *stx* subtyping provided a good level of sensitivity and specificity. In addition, the
52 WGS data also showed there were a significant proportion of strains harbouring multiple alleles
53 of the same Stx subtype associated with clinical disease in England.

54 **Introduction**

55 Shiga toxin-producing *Escherichia coli* (STEC) are a rare but potentially fatal cause of
56 gastroenteritis. They are associated with a wide spectrum of disease ranging from mild to bloody
57 diarrhoea, through to haemorrhagic colitis and haemolytic uraemic syndrome (HUS) (1). The
58 main reservoir of STEC in England is cattle, although it is carried by other animals, mainly
59 ruminants. Transmission to humans occurs through direct or indirect contact with animals or their
60 environments; consumption of contaminated food or water, and through person-to-person
61 contact. Each year, there are approximately 900 cases of STEC O157 in England confirmed by
62 the Gastrointestinal Bacteria Reference Unit (GBRU) at Public Health England.

63 The primary STEC virulence factor responsible for the most serious outcomes of human
64 infection is Shiga toxin (Stx), an AB₅ toxin that targets cells expressing the glycolipid
65 globotriaosylceramide (Gb3), disrupting host protein synthesis and causing apoptotic cell death
66 (2). Renal epithelial cell membranes are enriched for Gb3 resulting in the kidneys bearing the
67 brunt of Stx toxicity and, in 5-10% of cases, this leads to the development of Hemolytic Uremic
68 Syndrome (HUS) (1). There are two types of Stx; Stx1 and Stx2 and both have multiple
69 subtypes. These subtypes can be differentiated using a PCR targeted at the encoding genes

70 described by (3). In addition, a web-based tool, VirulenceFinder has been developed which uses
71 a *de novo* assembly followed by BLAST approach to identify subtypes of Stx (4). This system
72 was shown to have good, but not perfect, agreement with PCR, although how it handles strains
73 that encode both *stx2a* and *stx2c* is uncertain as no strains that encoded both these subtypes
74 were examined (4). The ability of STEC to cause severe illness in humans is determined by
75 multiple bacterial factors (in addition to host factors), including Shiga toxin subtype. There is
76 evidence that the Stx2a subtype is significantly associated with progression to HUS (5, 6).

77 As part of a project investigating the utility of whole genome sequencing (WGS) for public health
78 surveillance and outbreak investigation of foodborne pathogens, high throughput, short read
79 Illumina GAll sequence data for 444 strains of STEC O157 isolated in England between 2009
80 and 2013 was obtained. We determined the presence, or absence of the Stx encoding genes
81 *stx1* and/or *stx2* genes in all 444 isolates of STEC O157 from the genome sequence data. Given
82 the link between Stx subtype and disease severity, we also sought to identify the *stx* subtypes
83 present using bioinformatics methods and to compare the results with those obtained using the
84 PCR scheme of (3).

85 WGS high throughput short read technologies are rapid and low cost compared to Sanger
86 sequencing but it was recognised early on that assembling short reads would be problematic (7).
87 A major difficulty in assembly is the presence of repeat sequences that are longer than the read
88 length. Furthermore the study by (3) clearly demonstrated that as well as a high level of similarity
89 between *stx2a*, *stx2c* and *stx2d* there is also considerable diversity within each of these
90 subtypes. The assembly of *stx* into one contig in strains of STEC O157 containing both *stx2a*
91 and *stx2c* is difficult because the regions of variation between these subtypes are concentrated
92 at the 5' and 3' ends of the coding DNA sequence (CDS), with a largely homogenous region in
93 the centre. Existing methods for subtyping *stx* from short read data have not been tested against
94 strains encoding *stx2a* and *stx2c* (4). This region of 100% identity is often longer than the typical
95 read length of short read sequencing technologies, so contiguous assembly of both subtypes
96 relies on information from the paired end reads, which has a limited ability to resolve repeats up
97 to the average fragment size (550-700 bp for Illumina Nextera mate-pair). The STEC O157
98 Sakai reference genome encodes 18 pro-phage that show a large degree of modularity and
99 similarity (8), this further complicates assembly of these regions (9). These difficulties have led
100 to a relative paucity of data on the presence of subtypes of Stx within the *E. coli* population
101 despite large WGS projects.

102 In this study a dual bioinformatic approach was taken, using both mapping and *de novo*

103 assembly to determine *stx* subtype. The results of the bioinformatic analysis were compared to
104 the results from the PCR typing method (3). In addition, the diversity within and between the *stx*
105 subtype genes were investigated and evidence that certain strains contained multiple copies of
106 the same *stx* subtype was assessed.

107 **Methods**

108 **Strain selection**

109 A total of 444 isolates of STEC O157 submitted to GBRU for confirmation and typing were
110 selected for sequencing, 365 from 2012 representing approximately one third of the culture
111 positive isolates (1002 total isolates) received by the reference laboratory that calendar year
112 from laboratories in England, Wales and Northern Ireland, 67 English historical isolates
113 submitted to GBRU between 1990 and 2011 and 12 isolates from 2013. The collection contained
114 strains from sporadic cases, known outbreaks, household clusters, and serial strains isolated
115 from the same patient. However, only sporadic strains and a single strain from any related cases
116 (e.g. household, outbreak) were included in the diversity and multiple allele analysis. A total of 18
117 phage types were represented.

118 **Sequencing**

119 Genomic DNA was fragmented and tagged for multiplexing with Nextera XT DNA Sample
120 Preparation Kits (Illumina) and sequenced at the Animal Health Veterinary Laboratory Agency
121 (Weybridge) using the Illumina GAII platform with 2x150bp reads. Multiplexing allowed 96
122 samples to be sequenced per run. Sequencing data with a phred score below 30 or a read
123 length below 50 were removed from the data set using Trimmomatic (11). FASTQ data is
124 available from the NCBI Short Read Archive, BioProject accession PRJNA248064.

125 **Subtyping of *stx* by assembly**

126 High quality reads were assembled using Velvet v1.2.03 (10) with k-mer chosen using VelvetK
127 (<http://bioinformatics.net.au/software.velvetk.shtml>). The resulting contigs were then compared
128 against a set of *stx* reference genes (*stx1a*, L04539.1; *1c*, Z36901.1; *1d*, AY170851.1; *2a*,
129 X07865.1; *2b*, X65949.1; *2c*, AB071845.1; *2d*, AY095209.1; *2e*, AJ249351.2; *2f*, AB472687.1;
130 *2g*, AY286000.1) using BLASTn within the BioPython framework (12). Only matches with an E-
131 value less than 1×10^{-20} were included in further analysis. For each strain, the length of the best-
132 matched sequence (in terms of the BLAST score) between the contigs and each *stx* reference
133 gene was calculated. For example where both *stx2a* and *stx2c* were present, there may be five
134 query sequences each of 600 bp. If three of them matched *stx2a* with the highest BLAST score,
135 and two of them matched *stx2c* with the best BLAST score, then *stx2a* would score 1200 and

136 *stx2c* would score 800.

137 **Subtyping of *stx* by mapping**

138 An alignment of *stx1a*, *stx1c*, *1d*, *2a*, *2b*, *2c*, *2d*, *2e*, *2f* and *2g* sequences (taken from Scheutz et
139 al., 2012) was generated using ClustalW within the MEGA 5 software package (13). Three bases
140 for each reference subtype that, when combined, had 100% sensitivity and specificity for each
141 subtype were identified. High quality sequencing reads were mapped to a set of reference *stx*
142 genes (same genes as BLAST approach described above) using BWA-MEM ([http://bio-
143 bwa.sourceforge.net/](http://bio-bwa.sourceforge.net/)). Reads that mapped to more than one place in the reference set (i.e.
144 ambiguous reads) were removed from the resultant SAM file using Samtools (14). If at least 10
145 reads and 90% of the total reads concordantly mapped to all three discriminatory positions for a
146 specific subtype, then a positive match was returned for that subtype.

147 **Determination of the presence of multiple alleles of the same *stx* subtype by mapping 148 depth**

149 Multiple copies of the same *stx* allele could be identified using two complementary approaches.
150 In the first approach, reads were mapped to the *stx* reference genes, with ambiguous mapping
151 allowed. Then the coverage of each *stx* allele, which had been identified by the mapping and
152 assembly methods described above, was calculated using the Samtools 'depth' option. A
153 distribution of mapping depth in all the strains that were positive for one particular *Stx* subtype
154 was plotted revealing a bimodal distribution with the higher mode approximately twice the lower
155 mode. The lower mode represented strains with only one copy of *stx* and the higher mode
156 represented strains with multiple alleles of *stx*. There was no bimodal distribution of mapping
157 depth for strains that encoded both *stx2a* and *stx2c*, due to the redundant mapping between
158 these two strains. For example, if a strain encoded *stx2a* only and mapped to an *stx2c* reference
159 gene, it showed approximately one third of the average coverage compared to if it were mapped
160 to an *stx2a* gene. This cross-mapping meant that multiple alleles of the same *stx* subtype could
161 not be detected in strains that encoded both *stx2a* and *stx2c*.

162 In the second approach, the bam file resulting from the mapping of the reads to the *stx* reference
163 set was parsed for mixed positions with the minority variant present in at least 25% of reads i.e.
164 one position in the reference gene was mapped by two different bases. Only strains that were
165 known to encode only one of *stx2a* or *stx2c* from the subtyping results were analysed, as the
166 high similarity between *stx2a* and *stx2c* can result in pseudo-mixed bases when compared with
167 *stx* reference genes. If there were mixed bases present in an alignment (where the depth was
168 greater than 20x and minority variant present in greater than 15% of reads), from a strain

169 encoding only one of *stx2a* or *stx2c*, the presence of multiple alleles of a specific *stx* subtype
170 that vary by at least one base was assumed to be present (supplementary materials Figure 1).

171 **Diversity of *stx* associated with STEC O157 in the England, Wales and Northern Ireland**

172 The *stx* genes that were successfully assembled into a single contig were extracted from the *de*
173 *novo* genome assemblies using BLAST and aligned. Only strains that subtyping had shown to
174 encode one of *stx2a* or *stx2c* were included in this part of the study. Where the complete
175 sequence of both *stxA* and *stxB*, including the intergenic region, was assembled into a single
176 contig, the CDSs were aligned and represented in minimum spanning trees generated using
177 Bionumerics v6 (<http://www.applied-maths.com/bionumerics>). Strains where the *stxA* and *stxB*
178 subunits could not be assembled into a single contig (e.g. due to the presence of multiple copies
179 of the same *stx* subtype with sequence variation between them), were aligned against a
180 reference gene and the resulting Sam file was parsed using custom python scripts to identify
181 variant positions. The sequences of *stx1a*, *stx2a* and *stx2c* present in the strains investigated
182 here were compared with a representative sample of *stx* subtype sequences in the National
183 Centre for Biotechnology Information (NCBI, <http://www.ncbi.nlm.nih.gov/>) nucleotide database to
184 assess diversity and identify novel alleles.

185 **Stx real-time qPCR and block-based subtyping PCR**

186 DNA was prepared by inoculating single colonies into 490ul distilled water and boiled in a water
187 bath for 10 mins. The real-time qPCR described by the European Union Reference Laboratory
188 (EURL) for *stx1* and *stx2* was performed as previously described (15). For the block-based
189 subtyping PCR, DNA was amplified on a block-based DNA Engine platform using the *stx*
190 subtyping primers and amplification parameters described by (3). Amplified DNA was
191 electrophoresed on a 2% gel, stained with ethidium bromide and visualised with UV light.

192 **Results**

193 **Stx subtyping of 444 STEC O157 in the UK – comparison of NGS and PCR**

194 Subtyping results from PCR and WGS were identical in 422/444 strains (Table 1), there was
195 agreement for 85 *stx2a* encoding strains, 153 *stx2a/stx2c* strains and 187 *stx2c* strains. When
196 the subtyping PCR was repeated for the 22 discordant strains, results for 442/444 strains were
197 identical. Of the two strains where PCR and sequencing were discordant, one strain was positive
198 for *stx2c* by PCR but no *stx2c* was identified in the sequencing data by the bioinformatics
199 algorithms described here, and one strain was positive for *stx2a* by sequencing that was not

200 detected by PCR. The strain that had a positive PCR result for *stx2c* but no corresponding result
201 in the WGS data had a very low level of mapping (54 reads, <7x average coverage) to *stx2c*.
202 This was not enough to definitively identify *stx2c* by either the mapping or assembly algorithms,
203 although is indicative of its presence. The *stx2a* gene sequence of the strain that was PCR
204 negative but that had *stx2a* reads identified in the WGS data, was analysed for mutations in the
205 primer binding sites, but none were identified.

206 **Detection of multiple alleles of *stx***

207 A subset of 349 sporadic strains of STEC O157 (i.e. not from same person, household or
208 outbreak) was investigated for the presence of multiple alleles of the same subtype of *stx*. The
209 detection of multiple copies of the same *stx* subtype was performed using two complimentary
210 methods (i) mapping and determining the short read coverage of a particular *stx* subtype relative
211 to the coverage of the whole genome and (ii) the detection of mixed bases (coverage >20x,
212 minority variant >15%, see supplementary material Figure 1) in an alignment to a single
213 reference gene.

214 The *stx1a* gene was detected, in 6 different combinations with other subtypes/alleles, in 100
215 strains from independent sources (Table 2). For clarity, the relative coverage of *stx1a* in three of
216 the observed combinations (totalling 77 strains) is presented (Figure 1). The relative coverage of
217 *stx1a* in all 6 combinations observed in the 100 *stx1a* strains can be seen in supplementary
218 material Figure 2. In Figure 1, a bimodal distribution was clear, with the higher mode being
219 approximately twice as high as the lower mode. There were 11 strains in the higher mode
220 (Figure 1). When the *stx1a* alignments were examined for the presence of mixed bases, there
221 were 97 strains with no mixed bases and three strains that had at least one mixed base position.
222 The relative coverage was examined in the context of the presence of mixed bases and strains
223 with no mixed bases had a median relative coverage of 1.7x whereas strains with at least one
224 mixed base had a mean coverage of 2.8x (Figure 1). There were nine strains without mixed
225 bases that had relative *stx1a* coverage closer to the average of mixed base position strains than
226 the average of no mixed bases suggesting that two identical copies of the *stx1a* gene were
227 present.

228 There were 210 isolates that encoded *stx2a*, either alone or in combination with other subtypes
229 (Table 2). For clarity, only the relative coverage of *stx2a* from the 73 strains that encoded only
230 *stx2a* were presented in Figure 2 (the relative coverage of *stx2a* in all strains which encoded this
231 subtype can be seen in supplementary material Figure 3). Inspection of the distribution of
232 coverage of the short reads in *stx2a* revealed at least two modes within the relative coverage of

233 *stx2a*, with the upper mode (1.8x) being twice that of the lower (0.9x) (Figure 2). Of the 70
234 strains that encoded *stx2a* but not *stx2c*, 31 (42%) had short read coverage in the upper mode
235 (1.8x), of which all 31 had mixed base positions in their alignments, indicating the presence of
236 two alleles of *stx2a*. When the mixed position data was compared with the relative coverage
237 distribution, the mean relative coverage of the strains with mixed positions of *stx2a* was 1.9x
238 while the coverage in strains with no mixed positions was 0.75x (Figure 2). There were 1 (n =
239 29), 2 (n = 1) or 3 (n = 1) positions with mixed bases between the alleles in the 31 strains with
240 multiple copies.

241 The relative coverage of the 279 isolates that encoded *stx2c* was calculated. For clarity, only the
242 relative coverage of the 139 strains that encoded *stx2c* but not *stx2a* are presented in Figure 3.
243 The relative coverage of *stx2c* in all 279 strains can be seen in the supplementary materials and
244 analysis of the distribution of relative *stx2c* coverage showed that the majority of these strains
245 fell into an approximately normal distribution around 1x relative coverage (Figure 3). Twelve
246 (8.6%) of the 139 strains had a relative coverage > 1.5x but no mixed base positions were
247 found.

248 **Diversity of *stx* associated with STEC O157 in the UK**

249 The diversity of *stx* found in a subset of 349 sporadic strains of STEC O157 (i.e. not from same
250 person, household or outbreak) was investigated. Ninety-seven complete *stx1a* genes from this
251 study were compared with nine *stx1a* alleles from NCBI, and a total of 16 variant positions were
252 identified along the 1392 bp length of the gene. Of the five different alleles present in the strains
253 investigated here, three were not present in the NCBI database (as of 06/23/14, Figure 4). The
254 most frequently observed allele accounted for 76 (78.3%) of the 97 assembled *stx1a* genes from
255 this study, while the second most frequently observed allele accounted for 16 (16.5%) *stx1a*
256 genes. Both the most frequently observed alleles had been previously identified in *E. coli*
257 O103:H2 (BAI33872.1) and *E. coli* O157:H7 (EF079675.1), respectively. The five remaining
258 *stx1a* genes comprised three different alleles, none of which had been previously submitted to
259 the NCBI database ((as of 06/23/14), although they were all within a single variant of previously
260 observed alleles (Figure 4).

261 The 38 fully assembled *stx2a* genes from this study were compared with 21 *stx2a* alleles from
262 the NCBI nucleotide database. There were a total of 48 variant positions in a 1442 bp alignment
263 of the 59 *stx2a* genes that included 25 different alleles (Figure 5). Of these 25 alleles, six were
264 present in the strains investigated here. The most frequently observed allele was a single variant
265 from a *stx2a* allele observed in *E. coli* O157 (AF524944.1), *E. coli* O111 and *E. coli* O145 and

266 was present in 18 (47.3%) of the strains in this study. The second most frequently observed
267 allele was present in 11 (28.5%) strains and was widely distributed, including in Bacteriophage
268 933W (X07865). The other nine strains represented four alleles, two of which had been identified
269 before. The remaining allele (from strain H124840173) was highly divergent from the other *stx2a*
270 alleles, with six SNPs compared to any previously identified *stx2a* gene and 11 variants
271 compared to the closest *stx2a* observed in this study. Interestingly, this strain was a sorbitol-
272 fermenting (SF) STEC O157, the only SF strains to be included in this study.

273 There were 132 fully assembled *stx2c* genes from this study that were compared with 18
274 previously identified *stx2c* alleles from NCBI. There was a total of 59 variant positions along the
275 1441 bp gene alignment of the 150 *stx2c* sequences, comprising 22 unique alleles, of which
276 seven were identified in the strains analysed here (Figure 6). The most frequently observed
277 allele accounted for 115 (87.1%) of the 132 fully assembled *stx2c* genes. This allele had been
278 previously observed in a single *E. coli* O157:H- strain (AB015057.1). The 17 other *stx2c* genes
279 represented six distinct alleles that, with one exception, were within two variants of the most
280 frequently observed allele (Figure 6). There were two strains encoding the most divergent
281 observed *stx2c* allele, with six variant positions compared with the most frequently observed
282 allele. This divergent allele had been previously identified in *E. coli* RM10648 (KF932369.1).

283 Although the complete gene sequence could not be determined for strains that had more than
284 one copy of a *stx* subtype, an alignment of the reads against a reference was analysed to
285 identify variant positions. Of the three strains with multiple alleles of *stx1a*, all three had the
286 same four variant positions. There was one SNP in all three multiple-*stx1a* strains that was not
287 previously identified in the *stx1a* sequences described above or in the NCBI reference
288 sequences. Of the 30 strains with multiple copies of *stx2a*, 28 had only a single variant position
289 that was the same in all 28 strains and that had been previously identified. Of the other two
290 strains, one had the same SNP as the 28 other mixed position strains and an additional SNP
291 that had not been previously observed in the strains described in this study above or in the NCBI
292 reference strains. The final strain had three unique mixed positions, all of which had been
293 previously observed in this study.

294 **Discussion**

295 In this study we have developed novel, robust and highly accurate methods for subtyping of *stx*
296 from short read sequence data, validating this method against PCR for 444 STEC O157 isolates.
297 Furthermore, we have mined the WGS data to show that a significant proportion of strains
298 encode multiple copies of the same subtype of Shiga-toxin gene. The diversity of *stx* genes from

299 STEC O157 in England was also elucidated.

300 There was over 95% initial agreement (422 of 444 strains) between WGS subtyping and PCR
301 subtyping in determining subtypes of *stx2* which shows that WGS is an acceptable method for
302 subtyping *stx* in O157. The strains where there were discrepancies between WGS and PCR
303 were subjected to a repeat subtyping PCR, after which all but two of the discrepancies became
304 concordant. One possible reason for the discrepancy between the initial and repeat PCR results
305 is the high stringency of the subtyping PCR. During a multi-centre evaluation of the subtyping
306 PCR, there were differences observed in the subtyping results obtained between different
307 laboratories and these were ascribed to the use of different reagents and thermocyclers, with the
308 main source of variability thought to be the use of different polymerase. While the *tag*
309 polymerase recommended by (3) was used here, variations in other laboratory reagents and
310 equipment may have resulted in the discrepancies. The excellent concordance between the
311 PCR and WGS results, even despite the problems associated with analysis of homologous
312 genes using short read data, provides evidence of the accuracy of the bioinformatics algorithm
313 showing that WGS could replace PCR for subtyping.

314 Using mapping coverage to detect multiple copies of the Stx phage has been described
315 previously using more challenging metagenome data (16). The novelty of this work is to use the
316 mapping coverage of *stx* relative to the average coverage of the whole genome to identify
317 strains encoding multiple alleles of the same *stx* subtype. There are *stx* sequences in the NCBI
318 database that indicate that multiple alleles of the same subtype encoded by the same strain
319 have been previously observed i.e. these sequences contain ambiguous bases. However, the
320 studies associated with these sequences make no mention of the possibility of multiple alleles
321 (17, 18, 19, 20). The presence of multiple alleles of the same *stx* type has been previously
322 identified by WGS (21), however this is the first study to present a large scale comparison of this
323 method with PCR subtyping. Some of the ambiguous positions in sequences in the NCBI
324 database were the same positions in *stx* as the mixed positions observed in this study,
325 supporting the evidence that multiple alleles exist and are present in the same strain. While
326 mapping of short reads has been successful at detecting multiple copies of the same subtype, it
327 has not been possible in strains that encode *stx2a* and *stx2c* due to ambiguous mapping
328 between these types (see supplementary materials, Figures 2-4). For characterisation of these
329 *stx2a/stx2c* strains, and full characterisation of the insertion sites and genomic context of the *stx*
330 alleles in strains encoding multiple copies of the same subtype, longer sequencing reads from
331 e.g. PacBio, are needed. There was also evidence that some strains of STEC O157 encoded
332 multiple alleles of both *stx1a* and *stx2c*, further characterisation of these strains to determine

333 whether they had a genetic determinant that made Stx phage acquisition more likely would be
334 interesting.

335 The functional implication of encoding multiple alleles of the same *stx* subtype remains unclear.
336 Three hypotheses to explain the 9% prevalence of strains encoding multiple alleles of the same
337 subtype are (i) these strains are more likely to cause symptomatic human disease (ii) these
338 strains have an fitness advantage which increases their chance of being present in the
339 environment (iii) carrying multiple alleles of the same subtype is 'merely' a side effect of the
340 recombinogenic capacity of Stx phage, which confers no phenotype. It is interesting that while
341 the multiple alleles of *stx2a* and *stx2c* always seem to have nucleotide differences, this is a
342 minority in the strains that encode multiple copies of *stx1a*. The close sequence relationship
343 between the multiple copies of the same subtypes raises the question whether they are derived
344 from multiple insertions by different Stx phage, or a phage or *stx* gene duplication.

345 This study reports on the diversity of *stx* observed in STEC O157 in the UK (except Scotland)
346 between 1990 and 2013, with a focus on 2012. Although there are 10 described subtypes of
347 *stx1* and *stx2* combined (3), in an examination of 444 strains covering a wide temporal spread
348 and range of phage types, only three subtypes (*stx2a*, *stx2c* and *stx1a*) were observed. Previous
349 studies examining strains from cattle and humans similarly found only *stx2a*, *stx2c* and *stx1a*
350 (22). The most diverse *stx* identified here was *stx2a*, followed by *stx1a* and then *stx2c* (Figures
351 4-6). The majority of *stx2c* were of a single genotype, and all the novel alleles identified were
352 within a single SNP of the majority genotype. This difference in diversity observed between
353 *stx2a* and *stx2c* is interesting considering that the background diversity of these two subtypes is
354 largely similar (3). Further studies in this laboratory aim to determine the phylogenetic context of
355 isolates encoding these two subtypes. This study also described 10 novel alleles of *stx*, with the
356 most diverse being an *stx2a* sequence 6 SNPs from any previously described *stx2a*. The fact
357 that this diverse *stx* was observed in a sorbitol fermenting strain indicates that there may be a
358 significant reservoir of *stx* diversity other serotypes of STEC. The majority of novel alleles had
359 sequences that were single nucleotide variants to previously described sequences.

360 This study is the first to describe *stx* subtyping by both PCR and WGS methods in a large
361 number of strains of STEC O157. Both the PCR and WGS approaches to *stx* subtyping provided
362 a good level of sensitivity and specificity. The WGS data also showed that a significant
363 proportion of strains of STEC O157 harbour multiple alleles of the same Stx subtype. The
364 functional significance of multiple alleles of the same subtype remains unclear, although this is
365 the subject of on going work. Furthermore, the WGS analysis highlighted 10 novel alleles of *stx*

366 identified in this study and enabled us to study the diversity of *stx* sequences in a population of
367 STEC O157 associated with clinical disease in England.

368 **Acknowledgements**

369 **This work was funded through the NIHR Scientific Research and Development Fund –**
370 **Project Number #108061. We would like to acknowledge the contribution of the PHE**
371 **Bioinformatics Unit for their technical assistance.**

372 **References**

- 373 1. **Pennington H.** 2010. *Escherichia coli* O157. *Lancet* **376**:1428–35.
- 374 2. **Ethelberg S, Olsen KEP, Scheutz F, Jensen C, Schiellerup P, Engberg J, Petersen**
375 **AM, Olesen B, Gerner-Smidt P, Mølbak K.** 2004. Virulence Factors for Hemolytic
376 Uremic Syndrome. *Emerg. Infect. Dis.* **10**:842–847.
- 377 3. **Scheutz F, Teel LD, Beutin L, Piérard D, Buvens G, Karch H, Mellmann A, Caprioli A,**
378 **Tozzoli R, Morabito S, Strockbine N a, Melton-Celsa AR, Sanchez M, Persson S,**
379 **O'Brien AD.** 2012. Multicenter evaluation of a sequence-based protocol for subtyping
380 Shiga toxins and standardizing Stx nomenclature. *J. Clin. Microbiol.* **50**:2951–63.
- 381 4. **Joensen KG, Scheutz F, Lund O, Hasman H, Kaas RS, Nielsen EM, Aarestrup FM.**
382 2014. Real-time whole-genome sequencing for routine typing, surveillance, and outbreak
383 detection of verotoxigenic *Escherichia coli*. *J. Clin. Microbiol.* **52**:1501–10.
- 384 5. **Persson S, Olsen KEP, Ethelberg S, Scheutz F.** 2007. Subtyping method for
385 *Escherichia coli* shiga toxin (verocytotoxin) 2 variants and correlations to clinical
386 manifestations. *J. Clin. Microbiol.* **45**:2020–4.
- 387 6. **Luna-Gierke RE, Griffin PM, Gould LH, Herman K, Bopp CA, Strockbine N, Mody**
388 **RK.** 2014. Outbreaks of non-O157 Shiga toxin-producing *Escherichia coli* infection: USA.
389 *Epidemiol. Infect.* doi:10.1017/S0950268813003233
- 390 7. **Chaisson M, Pevzner P, Tang H.** 2004. Fragment assembly with short reads.
391 *Bioinformatics* **20**:2067–74.
- 392 8. **Asadulghani M, Ogura Y, Ooka T, Itoh T, Sawaguchi A, Iguchi A, Nakayama K,**
393 **Hayashi T.** 2009. The defective prophage pool of *Escherichia coli* O157: prophage-
394 prophage interactions potentiate horizontal transfer of virulence determinants. *PLoS*
395 *Pathog.* doi:10.1371/journal.ppat.1000408
- 396 9. **Hayashi T, Makino K, Ohnishi M, Kurokawa K, Ishii K, Yokoyama K, Han C, Ohtsubo**
397 **E, Nakayama K, Murata T, Tanaka M, Tobe T, Iida T, Takami H, Honda T, Sasakawa C,**
398 **Ogasawara N, Yasunaga T, Kuhara S, Shiba T, Hattori M.** 2001. Complete Genome
399 Sequence of Enterohemorrhagic *Escherichia coli* O157:H7 and Genomic Comparison with
400 a Laboratory Strain K-12. *DNA Res.* **8**:11–22.
- 401 10. **Zerbino DR, Birney E.** 2008. Velvet: algorithms for de novo short read assembly using
402 de Bruijn graphs. *Genome Res.* **18**:821–9.
- 403 11. **Bolger AM, Lohse M, Usadel B.** 2014. Trimmomatic: a flexible trimmer for Illumina
404 sequence data. *Bioinformatics* doi:10.1093/bioinformatics/btu170
- 405 12. **Cock PJ a, Antao T, Chang JT, Chapman B a, Cox CJ, Dalke A, Friedberg I,**
406 **Hamelryck T, Kauff F, Wilczynski B, de Hoon MJL.** 2009. Biopython: freely available
407 Python tools for computational molecular biology and bioinformatics. *Bioinformatics*
408 **25**:1422–3.
- 409 13. **Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S.** 2011. MEGA5:
410 molecular evolutionary genetics analysis using maximum likelihood, evolutionary

- 411 distance, and maximum parsimony methods. *Mol. Biol. Evol.* **28**:2731–9.
- 412 14. **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,**
413 **Durbin R.** 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
414 **25**:2078–9.
- 415 15. **Jenkins C, Lawson AJ, Cheasty T, Willshaw G a.** 2012. Assessment of a real-time PCR
416 for the detection and characterization of verocytotoxigenic *Escherichia coli*. *J. Med.*
417 *Microbiol.* **61**:1082–5.
- 418 16. **Loman NJ, Constantinidou C, Christner M, Rohde H, Chan JZ-M, Quick J, Weir JC,**
419 **Quince C, Smith GP, Betley JR, Aepfelbacher M, Pallen MJ.** 2013. A culture-
420 independent sequence-based metagenomics approach to the investigation of an outbreak
421 of Shiga-toxigenic *Escherichia coli* O104:H4. *JAMA.* **309**:1502–10.
- 422 17. **De Baets L, van der Taelen I, De Filette M, Pie D, Allison L, De Greve H,**
423 **Hernalsteens J, Imberechts H.** 2004. Genetic Typing of Shiga Toxin 2 Variants of
424 *Escherichia coli* by PCR-Restriction Fragment Length Polymorphism Analysis. *Appl.*
425 *Environ. Microbiol.* **70**:6309–6314.
- 426 18. **Lee JE, Reed J, Shields MS, Spiegel KM, Farrell LD, Sheridan PP.** 2007. Phylogenetic
427 analysis of Shiga toxin 1 and Shiga toxin 2 genes associated with disease outbreaks.
428 *BMC Microbiol.* **7**:109.
- 429 19. **Asakura H, Makino S, Kobori H, Watarai M, Shirahata T, Ikeda T, Takeshi K.** 2001.
430 Phylogenetic diversity and similarity of active sites of Shiga toxin (Stx) in Shiga toxin-
431 producing *Escherichia coli* (STEC). *Epidemiol. Infect.* **127**:27–36.
- 432 20. **Hegde A, Ballal M, Shenoy S.** 2012. Detection of diarrheagenic *Escherichia coli* by
433 multiplex PCR. *Indian J. Med. Microbiol.* **30**:279–284.
- 434 21. **Eppinger M, Mammel MK, Leclerc JE, Ravel J, Cebula TA.** 2011. Genomic anatomy of
435 *Escherichia coli* O157: H7 outbreaks. *PNAS* **108**:20142–20147.
- 436 22. **Mellor GE, Besser TE, Davis M a, Beavis B, Jung W, Smith H V, Jennison A V, Doyle**
437 **CJ, Chandry PS, Gobius KS, Fegan N.** 2013. Multilocus genotype analysis of
438 *Escherichia coli* O157 isolates from Australia and the United States provides evidence of
439 geographic divergence. *Appl. Environ. Microbiol.* **79**:5050–8.

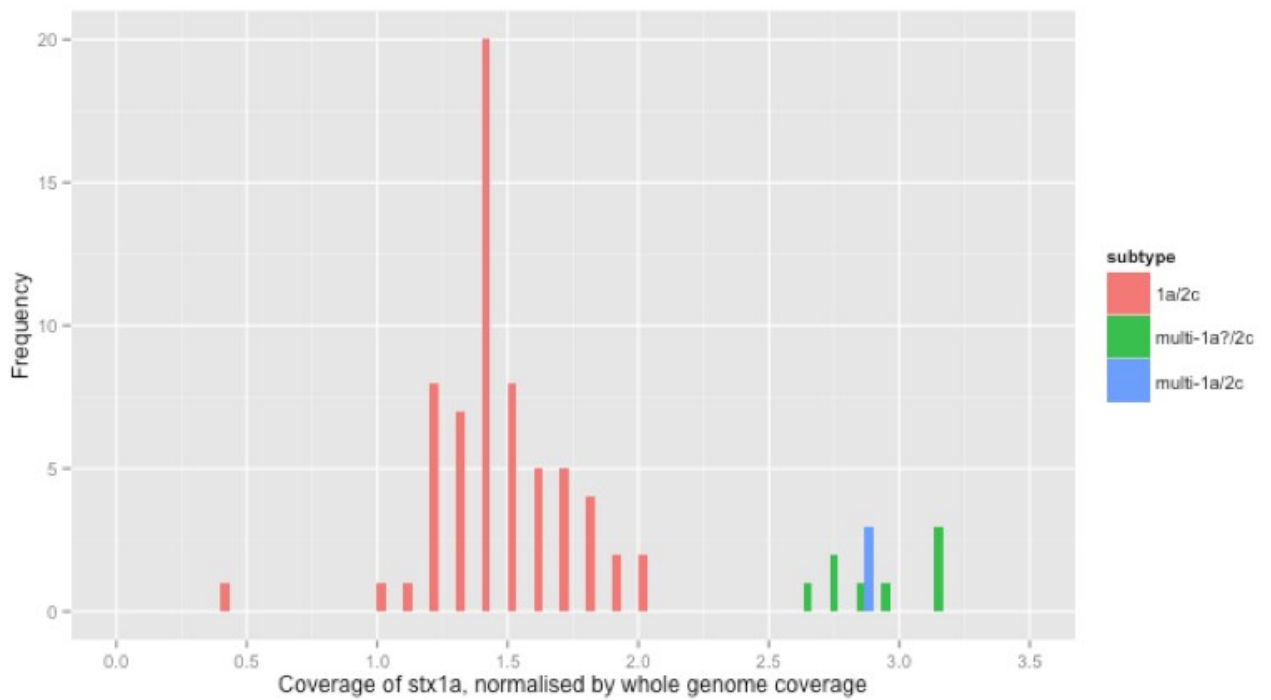
440 **Tables and Figures**

441 Table 1: Comparison of stx2 subtyping of 444 strains by sequencing and PCR. Strains that had
 442 discrepant results between sequencing and PCR were subjected to a 'second pass' PCR.

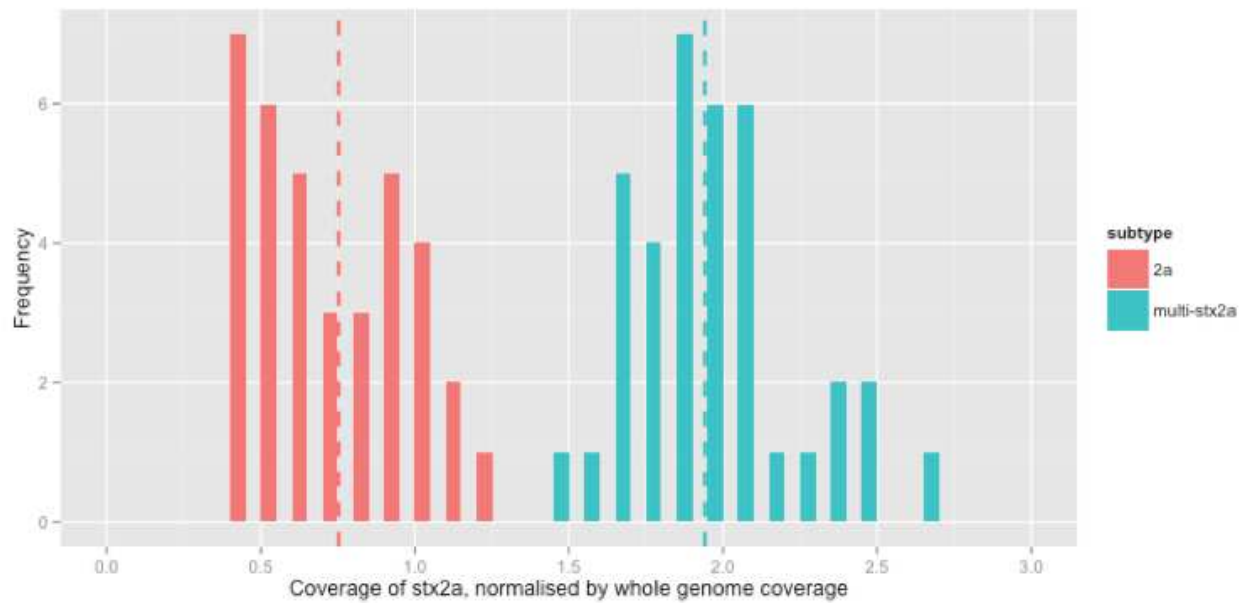
Subtype	Sequencing Results	Subtyping PCR results - 1st pass	Subtyping PCR results - 2nd pass
2a	82	89	82
2c	194	196	196
2a/2c	167	155	166
No result	1	4	0
Total	444	444	444

443 Table 2: Frequency of stx subtype profiles including stx1, derived from WGS analysis, not
 444 including outbreak strains. When a multi subtype result has a '?', it indicates that the only
 445 evidence suggesting the presence of multiple copies was the relative coverage (as opposed to
 446 having mixed positions as well).

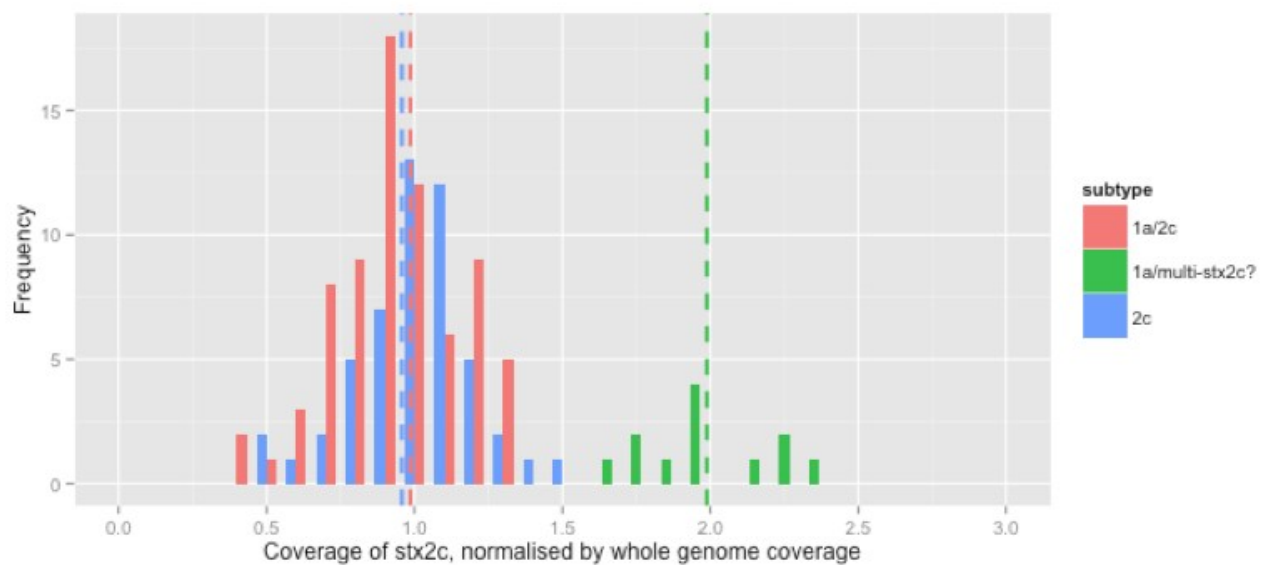
stx profile	Frequenc y
1a/2a	9
1a/2a/2c	3
1a/2c	64
1a/multi-stx2c?	10
2a	30
2a/2c	136
2c	51
multi-1a?/2c	9
multi-1a/2c	3
multi-stx2a	31
multi-stx2c?/multi-1a?	2
No stx detected	1



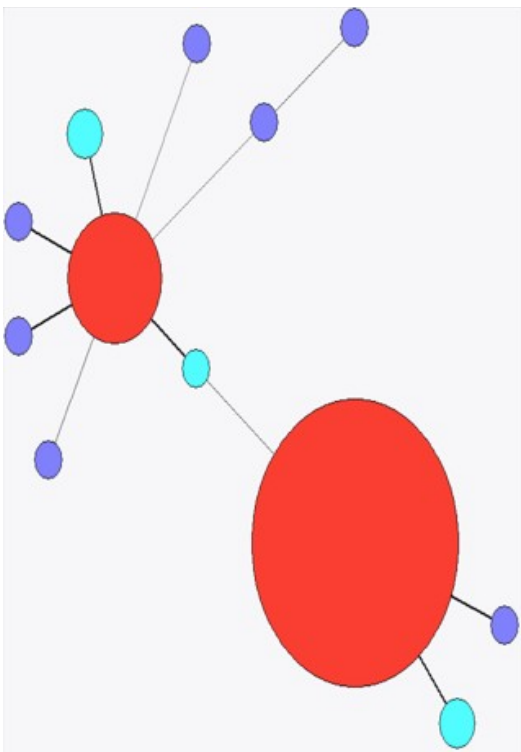
447 Figure 1: Histogram of coverage of stx1a normalised by whole genome coverage



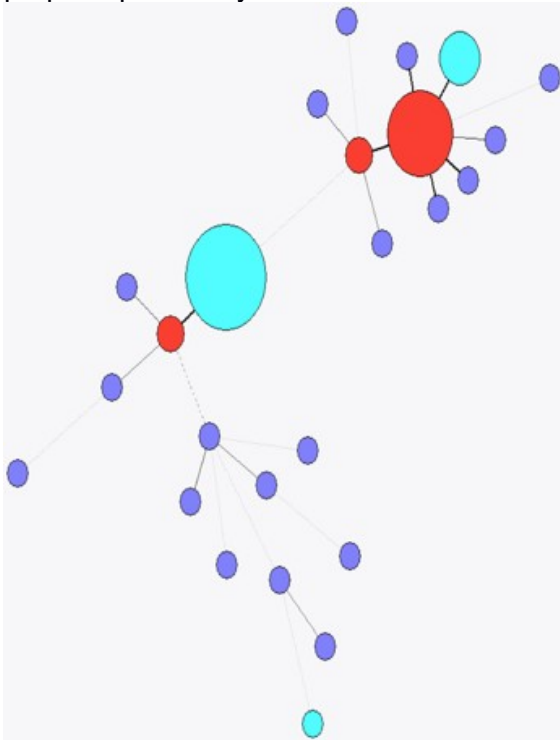
448 Figure 2: Histogram of coverage of stx2a normalised by whole genome coverage.



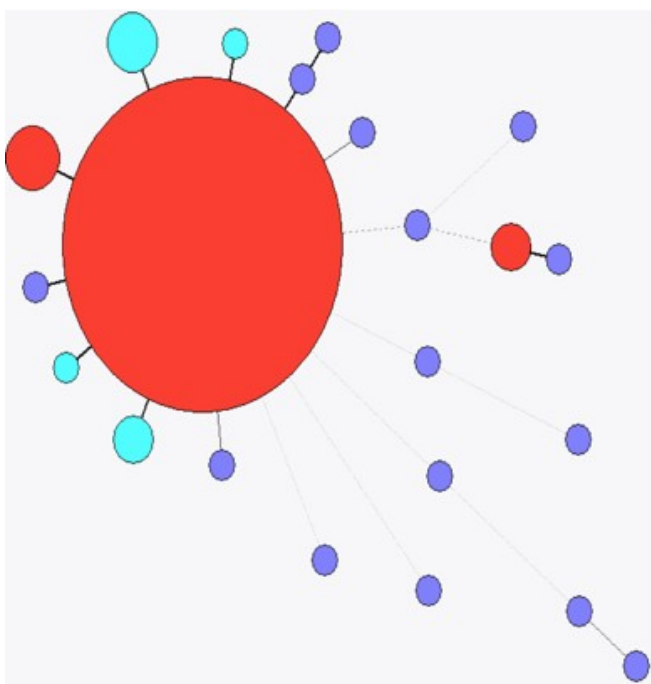
449 Figure 3: Histogram of coverage of stx2c normalised by whole genome coverage



450 Figure 4: Minimum spanning tree *stx1a*. Red = previously identified and observed in this study,
451 purple = previously identified but not observed in this study, light blue = novel allele.



452 Figure 5: minimum spanning tree *stx2a*. Colour as in Figure 4



453 Figure 6: Minimum spanning tree of *stx2c*. Colour as in Figure 4.

454 **Supplementary figures**

455 Supplementary Figure 1: Mixed position in an *stx2a* gene. Variant from reference highlighted.

456 Supplementary Figure 2: Histogram of coverage of *stx1a* normalised by whole genome coverage

457 Supplementary Figure 3: Histogram of coverage of *stx1a* normalised by whole genome coverage

458 Supplementary Figure 4: Histogram of coverage of *stx1a* normalised by whole genome coverage