

1 Title: Genomic evidence for revising the *Escherichia* genus and description of *Escherichia*
2 *ruysiae* sp. nov.

3

4 Running title (51/54 characters incl. spaces): Genomic evidence for revising the *Escherichia*
5 genus

6

7 Boas C.L. van der Putten^{a,b,#}, Sébastien Matamoros^a, COMBAT consortium[†], Constance
8 Schultsz^{a,b}

9 ^aAmsterdam UMC, University of Amsterdam, Department of Medical Microbiology,
10 Meibergdreef 9, Amsterdam, Netherlands

11 ^bAmsterdam UMC, University of Amsterdam, Department of Global Health, Amsterdam
12 Institute for Global Health and Development, Meibergdreef 9, Amsterdam, Netherlands

13 [#]Corresponding author, email: boas.vanderputten@amsterdamumc.nl

14 [†]Members listed in the appendix.

15

16 **Abstract (223/250 words)**

17 The *Escherichia/Shigella* genus comprises four *Escherichia* species, four *Shigella* species, and
18 five lineages currently not assigned to any species, termed '*Escherichia* cryptic clades'
19 (numbered I, II, III, IV and VI). Correct identification of *Escherichia* cryptic clades is strongly

20 hindered by the indeterminate taxonomy of this genus. Furthermore, little is known about the
21 cryptic clades as reports of genomic data are scarce. Hence, we searched public databases for
22 whole-genome sequences of *Escherichia* cryptic clades and characterized these. Following a
23 genomic analysis of the *Escherichia/Shigella* genus, we also describe a new *Escherichia* species:
24 *Escherichia ruysiae* sp. nov (type strain OPT1704^T = NCCB 100732^T = NCTC 14359^T) and provide a
25 closed genome assembly based on Illumina and Oxford Nanopore Technologies sequencing. We
26 screened 79,911 Sequence Read Archive *Escherichia* records and detected 357 cryptic clade
27 strains (0.44%). Based on average nucleotide identity, these strains should be grouped in seven
28 distinct species: 1) *E. coli*, *Shigella* spp. and clade I; 2) Clade II; 3) *Escherichia ruysiae* sp. nov.
29 (formerly clades III and IV); 4) *E. marmotae* (formerly clade V); 5) Clade VI; 6) *E. albertii* and 7) *E.*
30 *fergusonii*. Notably, half of the clade I strains carried genes encoding shiga toxin, while ESBL-
31 and carbapenemase-encoding strains were also found.

32 In conclusion, we provide an improved overview of the *Escherichia/Shigella* genus and advance
33 our understanding of *Escherichia* cryptic clades.

34

35 **Importance (107/150 words)**

36 Correct definition and identification of bacterial species is essential for clinical and research
37 purposes. Groups of *Escherichia* strains - "*Escherichia* cryptic clades" - have not been assigned
38 to species, which causes misidentification of these strains. The significance of our research is
39 threefold. First, we detect 357 cryptic clade strains, many more than previously known. This can
40 serve as a resource for other researchers. Second, we show how these cryptic clades should be

41 assigned to existing or newly defined species. This could improve identification of cryptic clade
42 strains and *Escherichia* species. Finally, we characterize the genomes in detail, revealing
43 virulence genes encoded in the cryptic clade I genomes.

44

45 **Main text excluding references, legends, tables: 3476/5000 words**

46 Introduction

47 Within the *Escherichia* genus, four species are recognized; *E. coli*, *E. fergusonii*¹, *E. albertii*² and
48 most recently, *E. marmotae*³. *E. coli*, *E. fergusonii* and *E. albertii* have been associated with
49 animal and human disease⁴⁻⁶, while little is known about the clinical relevance and
50 characteristics of the most recently described *Escherichia* species *E. marmotae*. Some
51 *Escherichia* strains cannot be assigned to any of the four existing species. Based on genetic
52 data, these strains cluster into six groups, which were termed '*Escherichia* cryptic clades'^{7,8}.
53 Recently, cryptic clade V was formally recognized as a separate species (*E. marmotae*)³, leaving
54 five cryptic clades that have not been recognized as species. As these have not yet received
55 entries in taxonomic databases, isolates representing cryptic clades might be misidentified by
56 tools that identify bacterial species based on whole-genome sequence (WGS). Also, the
57 abundancy of these cryptic clade strains in public databases such as RefSeq (curated repository
58 of reference sequences) or the Sequence Read Archive (SRA, repository of raw sequence data)
59 is unknown.

60 Therefore, we have addressed two questions in this study:

- 61 1) How abundant are *Escherichia* cryptic clades in the RefSeq and Sequence Read Archive (SRA)
62 databases?
- 63 2) Can the *Escherichia* cryptic clades be assigned to existing or new species based on genomic
64 analysis?

65 Finally, we describe a novel *Escherichia* species, *Escherichia ruysiae* sp. nov., encompassing the
66 former cryptic clades III and IV.

67

68 Results

69 Assessing abundance of cryptic clade strains in RefSeq

70 We downloaded all 10,824 available *Escherichia* genomes from the RefSeq database and
71 identified 92 cryptic clade genomes using ClermonTyper⁹. Since ClermonTyper is only able to
72 detect strains belonging to known cryptic clades, we constructed a neighbor-joining phylogeny
73 of the 10,824 RefSeq *Escherichia* strains using mash¹⁰ and rapidNJ¹¹ to detect cryptic clade
74 strains missed by ClermonTyper. This analysis revealed two additional cryptic clade strains,
75 bringing the total of cryptic clade strains in RefSeq to 94, or 0.9% of all *Escherichia* genomes in
76 RefSeq.

77 We then screened 79,911 Sequence Read Archive (SRA) *Escherichia* entries for presence of
78 cryptic clade genomes. In short, we performed a rough assembly of the SRA data using SKESA
79 and used fastANI to detect cryptic clade strains, using RefSeq cryptic clade genomes as
80 references. 757 putative cryptic clade genomes were included in a neighborhood-joining
81 phylogeny and were checked using ClermonTyper. This led to the discovery of 357 cryptic clade
82 genomes in the SRA, or 0.45% of SRA *Escherichia* assemblies. We evaluated SKESA assemblies
83 using Quast¹², and excluded 13 assemblies due to inadequate assembly qualities or aberrant
84 assembly length. Since there are likely duplicates between SRA and RefSeq cryptic clade

85 genomes, we chose to continue with the 344 cryptic clade genomes from the SRA. An overview
86 of the process is provided in figure 1.

87

88 **Genomic characterization of cryptic clade genomes from SRA**

89 To visualize phylogenetic relationships between the closely related *Escherichia* and *Shigella*
90 genera, we included – in addition to the 344 cryptic clade genomes - 50 *E. albertii* genomes, 10
91 *E. fergusonii* genomes, 72 *E. coli* genomes from the ECOR collection¹³ and 13 *Shigella spp.*
92 genomes, all from RefSeq. We constructed a maximum-likelihood phylogeny using IQ-Tree¹⁴
93 from a SNP alignment produced by kSNP3¹⁵ (fig. 2).

94

95 Based on species identification using Kraken2¹⁶ all cryptic clade I, II, III, IV and VI strains were
96 classified as *E. coli*. 82 out of 92 *E. marmotae* (formerly cryptic clade V) strains were identified
97 correctly while the remaining *E. marmotae* strains were identified as *E. coli*.

98

99 **Resistance and virulence genes in cryptic clade genomes**

100 Genes encoding both subunits of shiga toxin were detected in 108/209 (52%) of cryptic clade I
101 strains. In 26 cryptic clade I strains, *stx1a* and *stx1b* were detected together, in 75 strains *stx2a*
102 and *stx2b* were detected together. In 7 clade I strains only one shiga toxin-encoding gene was
103 detected. Extended-spectrum beta-lactamase (ESBL) genes were detected in 8.1% and 2.2% of

104 cryptic clade I and cryptic clade V/*E. marmotae* genomes, respectively. One clade I strain and
105 two *E. marmotae*/clade V strains harbored carbapenemase-encoding genes.

106

107 **Assignment of cryptic clades to existing or new species**

108 Based on a 95% average nucleotide identity (ANI)^{17,18} 489 of the included *Escherichia* and
109 *Shigella* strains should be assigned to seven discrete species with no relationships between
110 these groups (table 1 and fig. 3). *E. coli* and cryptic clade I showed mean ANI values of 96.0%,
111 indicating *E. coli* and cryptic clade I should be assigned into a single species, possibly separated
112 into subspecies. Genomically, *Shigella spp.* belong to *E. coli* as a mean ANI is 97.6% with *E. coli*.
113 This is in line with earlier reports¹⁹.

114

115 Cryptic clade III and IV should be assigned to a single species, as the mean ANI values between
116 strains from these groups was 96.6%. Based on ANI analysis cryptic clade II, cryptic clade V/*E.*
117 *marmotae*, cryptic clade VI, *E. fergusonii* and *E. albertii* should all be assigned to separate
118 species.

119

120 **Table 1.** Mean ANI values between groups calculated using fastANI, for 489 included
121 *Escherichia* and *Shigella* strains. Values higher than the threshold for species delineation (95%)
122 are in bold.

	Clade I	Clade II	Clade III	Clade IV	Clade V/ <i>E. marmotae</i>	Clade VI	<i>E. albertii</i>	<i>E. coli</i>	<i>E. fergusonii</i>	<i>Shigella</i>
Clade I	98.5	91.6	92.3	92.5	91.0	94.0	89.9	96.0	92.1	95.8
Clade II		99.0	92.1	92.0	91.3	90.9	89.7	91.8	89.2	91.7
Clade III			98.8	96.6	92.3	91.6	89.8	92.6	89.4	92.0
Clade IV				99.0	92.2	91.7	89.9	92.8	89.5	92.9
Clade V/ <i>E. marmotae</i>					99.3	90.4	89.4	91.2	88.4	91.1
Clade VI						99.9	89.6	94.4	90.1	94.0
<i>E. albertii</i>							98.6	90.1	88.2	90.0
<i>E. coli</i>								97.7	91.1	97.0
<i>E. fergusonii</i>									98.7	91.0
<i>Shigella spp.</i>										98.3

123

124

125 **Identification and description of *Escherichia ruysiae* sp. nov.**

126 We discovered a cryptic clade IV strain in our collection, previously identified as *E. coli* in the
 127 COMBAT study, which investigated acquisition of ESBL-producing Enterobacteriaceae (ESBL-E)
 128 during international travel²⁰. We characterized this isolate, OPT1704^T, in detail.

129 The strain was isolated from a faecal sample provided immediately after return from a one-
 130 month journey to several countries in Asia. No ESBL-E were detected in a faecal sample
 131 provided directly before departure, suggesting the ESBL gene, and possibly strain OPT1704^T,
 132 were acquired during travel. The traveller reported diarrhoea during travel and did not report
 133 antibiotic usage during travel. No ESBL-E were isolated in follow-up samples, suggesting loss of
 134 the OPT1704^T strain or the ESBL gene within one month after return from travel.

135

136 Strain OPT1704^T formed circular, grey-white colonies on a COS sheep blood agar plate.
137 Individual cells were observed under a light microscope and were rod-shaped and
138 approximately 1 by 2 μm in size. The strain was shown to be Gram-negative, non-motile,
139 oxidase-negative and catalase-positive. The strain was capable to grow in the absence of
140 oxygen. On COS blood plates, it showed growth in the temperature range of 20-42 °C, and no
141 growth at 4 °C or at 50 °C or higher. The strain was also able to grow in NaCl concentrations
142 ranging from 0% to 6% in liquid lysogeny broth. MALDI-TOF (Bruker) and Vitek2 (BioMérieux)
143 systems both identified OPT1704^T as *E. coli* with high confidence scores (score>2 for MALDI-TOF
144 and “Excellent identification” for Vitek2).

145

146 **Whole-genome sequence analysis**

147 The whole-genome DNA sequence of strain OPT1704^T was determined using Illumina HiSeq and
148 Oxford Nanopore Technologies (ONT) sequencing platforms. The Illumina sequencing run
149 yielded a total of 6.3×10⁶ paired-end reads, with a mean read length of 151 bp. Illumina reads
150 were downsampled using seqtk (version 1.3-r106, <https://github.com/lh3/seqtk>) to provide a
151 theoretical depth of coverage of 100X with the assumption that the OPT1704^T has a genome
152 size of approximately 5×10⁶ bp. The ONT sequencing run yielded a total of 2.5×10⁴ reads, with a
153 mean read length of 9078 bp before filtering. ONT reads were filtered on length and on read
154 identity using Illumina reads, leaving 1.5×10⁴ reads with a mean length of 12580 bp. This
155 provided a theoretical depth of coverage of ~38X of ONT reads. Assembly using both Illumina

156 and Nanopore reads resulted in a complete genome, consisting of one circular chromosome
157 and one circular plasmid. GC content of the OPT1704^T genome was 50.6%.

158
159 Putative resistance and virulence genes were predicted from the draft genome using ABRicate
160 with the CARD²¹ and VFDB²² databases. OPT1704^T harbours 6 resistance genes that are typically
161 plasmid-mediated in *E. coli*, associated with reduced susceptibility to fluoroquinolones (*qnrS1*),
162 aminoglycosides (*aph(6)-Id* & *aph(3'')-Ib*), cephalosporins (*blaCTX-M-14*), trimethoprim (*dfrA14*)
163 and sulphonamides (*su2*). This is in line with the reduced susceptibility to fluoroquinolones
164 (norfloxacin, MIC: 2 mg/L and ciprofloxacin, MIC: 0.5 mg/L by Vitek2), cephalosporins
165 (cefuroxime, MIC: >32 mg/L and cefotaxime, MIC: 4 mg/L) and trimethoprim-sulfamethoxazole
166 (MIC: >8 mg/L). Furthermore, several putative virulence genes were predicted from the draft
167 genome sequence associated with siderophore function (*chuX*, *entS*, *fepABD*), fimbriae
168 (*fimBCDGI*), a type II secretion system (*gspGHI*) and capsular polysaccharide biogenesis (*kpsD*).
169 The diarrhoeal symptoms of the traveller contributing this strain would not be expected based
170 on these predicted virulence genes.

171
172 Next, we calculated 16S rRNA sequence similarities, ANI values and digital DNA:DNA
173 hybridisation (dDDH) values between OPT1704^T and type strains of the four other *Escherichia*
174 species. This time, we used three separate tools to calculate ANI (fastANI, OrthoANIu and ANI
175 calculator from Enveomics)²³⁻²⁵, since scalability was no concern. 16S rRNA sequence
176 similarities did not completely warrant assignment of OPT1704^T to a novel species, but ANI

177 analysis and dDDH did support assignment of OPT1704^T to a novel species (table 2). This novel
178 species was assigned *E. ruysiae sp. nov.* with OPT1704^T as the proposed type strain. As our
179 earlier analysis shows cryptic clade III clusters with cryptic clade IV on the species level, we
180 propose that *E. ruysiae sp. nov.* encompasses both cryptic clade III and IV.

181

182 **Table 2.** Comparison of OPT1704^T whole-genome sequence with type strains of *E. albertii*, *E.*
183 *coli*, *E. fergusonii* and *E. marmotae*. In bold are the values that suggest assignment of OPT1704^T
184 to a novel species (<98.9% 16S rRNA sequence similarity, <95% ANI, <70% dDDH).

	<i>E. ruysiae sp. nov.</i> OPT1704 ^T				
	16S rRNA sequence similarity (%)	ANI (%, fastANI)	ANI (%, OrthoANIu)	ANI (% ANI calculator Enveomics)	dDDH (%)
<i>E. albertii</i> NBRC107761 ^T	98.5	90.0	90.0	89.2	39.8
<i>E. coli</i> ATCC11775 ^T	98.5	92.8	92.4	92.0	48.3
<i>E. fergusonii</i> ATCC35469 ^T	99.4	89.4	88.2	89.7	36.7
<i>E. marmotae</i> DSM 28771 ^T	99.4	92.2	92.2	91.4	47.1

185

186 Discussion

187 We propose an updated taxonomy of the *Escherichia* genus that includes seven *Escherichia*
188 species: 1) *E. coli*, cryptic clade I and *Shigella spp.*; 2) cryptic clade II; 3) *E. ruysiae sp. nov.*
189 (formerly cryptic clades III and IV); 4) *E. marmotae* (formerly cryptic clade V); 5) cryptic clade VI;
190 6) *E. fergusonii* and 7) *E. albertii*.

191 Contrary to the present study and possibly because a smaller dataset was used, Luo et al.
192 (2011) observed a 'genetic continuum' between *E. coli* and cryptic clades. In line with findings
193 by Walk (2015) we did not find strains that share an ANI > 95% between the seven groups,
194 indicating discrete grouping²⁶. Although they should be assigned to a single species based on
195 ANI analysis, a clear separation between *E. coli/Shigella spp.* and clade I is observed, possibly at
196 the subspecies level. The same holds for cryptic clades III and IV. For the other cryptic clades, no
197 significant separation in population structure is visible, indicating these clades are genetically
198 more homogeneous. The species clusters we find are also found by the Genome Taxonomy
199 Database (GTDB), where cryptic clade II is *Escherichia sp. 001660175*, cryptic clade VI is
200 *Escherichia sp. 002965065* and *E. ruysiae sp. nov.* is *Escherichia sp. 000208585*²⁷.

201
202 Throughout our comprehensive scanning of RefSeq and SRA, we attempted to maximize the
203 diversity of *Escherichia* genomes captured. We built the phylogeny of >10,000 RefSeq genomes
204 which included all *Escherichia* cryptic clades described in the literature. As we based our SRA
205 search on the RefSeq cryptic clade genomes, we cannot exclude that we missed novel cryptic
206 clades present in the SRA but not in RefSeq. Additionally, we had to raise the threshold with

207 which we screened the SRA for cryptic clade I strains. Possibly, strains that are genetically
208 different could have been missed by choosing this strategy.

209

210 Based on genomic analyses of the *Escherichia* genus and characterization of strain OPT1704^T,
211 we proposed a novel species, *Escherichia ruysiae* sp. nov. Analysis of 16S rRNA genes of *E.*
212 *ruysiae* and type strains of other *Escherichia* spp. did not warrant assignment to a novel species,
213 while ANI and dDDH analysis of the same genomes did. The International Journal of Systematic
214 and Evolutionary Microbiology follows results of ANI or dDDH analysis instead of 16S rRNA
215 analysis when discrepancies are observed, meaning strain OPT1704^T should be assigned to a
216 novel species¹⁷. Assignment of cryptic clades III and IV as a novel species *E. ruysiae* means
217 inclusion of the species in the NCBI taxonomy database, as happened with *E. marmotae*/cryptic
218 clade V before. As the NCBI taxonomy serves as the reference for many genomic classification
219 tools, assigning cryptic clades III and IV to a novel species is the most effective way to stimulate
220 correct classification of these strains.

221 Additionally, the assignment of *E. ruysiae* sp. nov. could be useful for recent efforts of NCBI to
222 retrospectively correct RefSeq entries using ANI, if the entries are assigned to the wrong species
223 ([https://ncbiinsights.ncbi.nlm.nih.gov/2019/02/06/correct-existing-taxonomic-info-genbank-ani-](https://ncbiinsights.ncbi.nlm.nih.gov/2019/02/06/correct-existing-taxonomic-info-genbank-ani-analysis/)
224 [analysis/](https://ncbiinsights.ncbi.nlm.nih.gov/2019/02/06/correct-existing-taxonomic-info-genbank-ani-analysis/)). However, this effort is also based on the NCBI taxonomy, which is why it is important
225 that cryptic clades that should be assigned as separate species are formally recognized as such.

226

227 All standard identification methods we used in this study – genomic and phenotypic –
228 concluded strain OPT1704^T is an *E. coli* strain. Genomic identification methods such as Kraken2
229 often use the RefSeq database in combination with the NCBI taxonomy to classify strains.
230 Cryptic clade strains are present in the RefSeq database, but are always labeled as either
231 indeterminate *Escherichia species* or as *E. coli*, as cryptic clades are missing from the NCBI
232 taxonomy. This causes Kraken2 to misidentify cryptic clade strains. Adding entries for cryptic
233 clades in the NCBI taxonomy (e.g. as formally recognized species) should solve this issue. It is
234 not unimaginable that the same issue holds for phenotypic characterization methods such as
235 the MALDI-TOF or Vitek2 systems, meaning some of the reference strains in those databases
236 could be wrongly labeled. Another possible explanation for the Vitek2 results is that it is
237 challenging to differentiate cryptic clades from *E. coli* based on biochemical properties. Walk et
238 al. (2009)⁸ assessed 31 biochemical markers and attempted to differentiate cryptic clades from
239 *E. coli* based on the biochemical profiles, which was achieved only partially. However, high-
240 throughput metabolomic analyses of cryptic clade strains have not been performed yet to the
241 best of our knowledge, which means there could be differentiating biochemical markers waiting
242 to be discovered.

243

244 Cryptic clade I genomes in our collection harbored more virulence and resistance genes
245 compared to other clades. Notably, almost half of clade I genomes harbored the genes
246 necessary to produce either Shiga toxin 1 or 2. It might be that this is partly due to a sampling
247 bias – toxigenic clade I strains might be whole-genome sequenced more often – but genes
248 encoding Shiga toxin were not nearly as abundant in other cryptic clades. In fact, cryptic clade I

249 might represent an eighth non-O157 *stx+* *Escherichia* lineage²⁸. Our analysis definitely confirms
250 the earlier suspicion of Walk (2015) that a relatively high percentage cryptic clade I strains
251 harbor shiga toxin²⁶. One positive point out of the perspective of public health is that none of
252 the *stx+* clade I genomes harbored plasmid-mediated genes conferring resistance to extended-
253 spectrum beta-lactams, carbapenems or colistin.

254
255 Our study provides an improved and systematic taxonomy of the *Escherichia/Shigella* genus,
256 making effective use of public sequence databases. Based on this analysis, we describe a novel
257 *Escherichia* species, *E. ruysiae* sp. nov., on a phenotypic and genomic level. The genomic
258 analysis of *Escherichia* population structure and the description of *E. ruysiae* should aid the
259 correct identification of *Escherichia* species in the future.

260 Materials and Methods

261 Assessing abundance of cryptic clade strains in RefSeq and SRA

262 10,824 *Escherichia* genome assemblies were downloaded from the RefSeq database on July
263 17th, 2018 using ncbi-download (version 0.2.6, [https://github.com/kblin/ncbi-genome-](https://github.com/kblin/ncbi-genome-download)
264 [download](https://github.com/kblin/ncbi-genome-download)). The downloaded genomes were processed using ClermonTyper (version 1.3.0)⁹ to
265 detect putative cryptic clade strains. To check for false negatives in the ClermonTyper analysis,
266 a neighbor-joining tree was constructed including all 10,824 *Escherichia* assemblies using mash
267 (version 2.1)¹⁰, the square_mash script from PySEER (version 1.2.0)²⁹ and rapidNJ (version
268 2.3.2)¹¹.

269 79,911 *Escherichia* entries were assembled from the Sequence Read Archive (SRA) using SKESA
270 (version 2.3.0)³⁰, assembling with a single kmer size of 51 without using read pairing
271 information and without removing adapter sequences. FastANI (version 1.1)²³ was used to find
272 genomes similar to any of the RefSeq cryptic clade genomes among these 79,911 assemblies.
273 Thresholds were set at 95% for clades II to VI. All strains with an ANI above the threshold for a
274 particular clade were selected for further analyses. Similarity was expressed in %ANI above
275 certain thresholds. Since cryptic clade I strains are similar to *E. coli*, a higher threshold for clade
276 I was chosen as using a 95% ANI threshold resulted in >50,000 hits. The lowest ANI value
277 between any two RefSeq cryptic clade I strains was 97.6%, so we chose a threshold of 96.5% to
278 reduce the probability of false negatives as much as feasible.

279 All SRA strains that showed similarity to RefSeq cryptic clade strains were re-assembled using
280 SKESA with default settings. Finally, a mash distance-based phylogeny was produced of all
281 reassembled strains to confirm phylogenetic placement of strains.

282

283 **Genomic characterization of cryptic clade strains**

284 Assembly metrics of all SRA assemblies were analyzed using Quast (version 4.5)¹². Additionally,
285 all SRA assemblies were screened using ABRicate (version 0.8.10,
286 <https://github.com/tseemann/abricate>) for virulence genes (using the VFDB core database,
287 downloaded 31st October 2018)²² or antimicrobial resistance (using the CARD database,
288 downloaded 31st October 2018)²¹. Kraken2 (version 2.0.7-beta) was used with the 8GB

289 MiniKraken database to perform species identification of cryptic clade genomes. Genome-wide
290 ANI was calculated using fastANI (version 1.1)²³ with default settings.

291

292 **Phylogenetic analysis**

293 A total of 489 genomes (344 cryptic clade genomes and 145 *E. albertii*, *E. fergusonii*, *E. coli* and
294 *Shigella spp.*) were used to construct a maximum likelihood phylogeny. kSNP3 (version 3.1)¹⁵
295 was used to extract SNPs from kmers present in all genomes. This core SNP matrix from kmers
296 was used as input for IQ-tree (version 1.6.6)¹⁴, which calculated a maximum-likelihood
297 phylogeny under the GTR-GAMMA model, correcting for ascertainment bias using 1000
298 ultrafast bootstraps³¹.

299

300 **Average Nucleotide Identity (ANI) analysis of 489 *Escherichia* and *Shigella* strains**

301 ANI analysis was performed using fastANI (version 1.1) which efficiently and accurately calculates
302 ANI for a large genome collection. 489 *Escherichia* and *Shigella* strains were used as query and
303 reference, providing a total of 238,632 comparisons (without self-comparisons). Mean ANI
304 values were calculated between all strains from compared groups. Subsequently, fastANI
305 output was written to a full matrix using the square_mash script of PySEER²⁹. Values lower than
306 95% were removed in R (version 3.5.2) and a Cytoscape network was created using the
307 graph_from_adjacency_matrix function from the igraph package (version 1.2.4)³², the jsonlite
308 package (version 1.6)³³ and the toCytoscape function from

309 [rest-R](#). The network was plotted in Cytoscape (version 3.7.1)³⁴ and layout was provided through
310 a weighted Prefuse Force Directed Layout algorithm distributed with Cytoscape.

311

312 **Phenotypic characterization of *Escherichia ruysiae* sp. nov.**

313 *Escherichia ruysiae* sp. nov. strain OPT1704^T strain was grown on COS sheep blood agar plates
314 (BD) at 37 °C unless stated otherwise. Anaerobic growth was assessed by growing OPT1704^T in
315 a sealed container with a BD Anaerobic GasPak within. A hanging drop preparation was used to
316 detect motility. The VITEK 2 system using a GN ID card (BioMérieux) and the MALDI-TOF system
317 (Bruker) were used with default settings for phenotypic identification of OPT1704^T. The Vitek
318 system was also used to obtain antimicrobial susceptibility profiles for multiple antibiotic
319 classes. EUCAST v9.0 clinical breakpoints were used to classify MICs into S/I/R categories.

320

321 **Genomic characterization of *Escherichia ruysiae* sp. nov.**

322 DNA for Illumina whole-genome sequencing was extracted using the Qiagen Blood and Tissue
323 kit according to the manufacturer's instructions. The sequencing library was prepared using the
324 Kapa HTP Library Preparation kit and subsequently sequenced on an Illumina HiSeq 4000
325 platform. Adapter sequences were removed using fastp³⁵.

326 DNA for Oxford Nanopore Technologies whole-genome sequencing was extracted using the
327 Qiagen MagAttract HMW DNA kit according to the manufacturer's instructions. The sequencing
328 library was prepared using the Nanopore rapid barcoding kit (SQK-RBK004) and sequenced in a

329 multiplexed run. Basecalling was performed using MinKNOW in fast basecalling mode and qcat
330 (version 1.0.7, <https://github.com/nanoporetech/qcat>) was used to demultiplex reads. Reads
331 were filtered using Filtlong (version 0.2.0, <https://github.com/rrwick/Filtlong>) where 90% of
332 best reads were retained, judged on read length and read identity based on Illumina
333 sequencing data. Filtered Nanopore reads and subsampled Illumina reads were used to perform
334 hybrid assembly using Unicycler (version 0.4.6, Wick 2017).

335 ABRicate was used as described above to predict resistance and virulence genes. 16S rRNA
336 gene sequence was extracted using barnap (version 0.9,
337 <https://github.com/tseemann/barnap>). Digital DNA:DNA hybridization scores were calculated
338 using the DSMZ webtool³⁶ and ANI was calculated using OrthoANIu²⁵, Enveomics ANI
339 calculator²⁴ and fastANI²³.

340

341 Acknowledgments

342 The COMBAT study was funded by Netherlands Organization for Health, Research and
343 Development (ZonMw; 50-51700-98-120) and EU-H2020 programme (COMPARE, 643476). The
344 authors would like to thank Rob Weijts and Patricia Brinke for their help in phenotypic
345 characterization of type strain OPT1704^T of *Escherichia ruysiae* sp. nov.

346

347 Data availability

348 All supplementary data is available via FigShare. This includes a Readme
349 (<https://doi.org/10.6084/m9.figshare.9824633.v1>), a list of accession numbers for the 10,824
350 *Escherichia* genomes downloaded from NCBI (<https://doi.org/10.6084/m9.figshare.9824234>), a
351 list of accession numbers for the 79,911 *Escherichia* entries assembled from the SRA using
352 SKESA (<https://doi.org/10.6084/m9.figshare.9824261>), the re-assemblies of the 344 selected
353 cryptic clade strains (<https://doi.org/10.6084/m9.figshare.9824270>) and a summary of the
354 bioinformatic commands used (<https://doi.org/10.6084/m9.figshare.9777746>). Illumina and
355 ONT fastq files are available on ENA under accession numbers ERR3518913 and ERR3518914,
356 respectively. Illumina and ONT sequencing data, as well as Unicycler assembly of strain
357 OPT1704^T are available on ENA under the project number PRJEB34275. Pure cultures of strain
358 OPT1704^T are available through the Netherlands Culture Collection of Bacteria of the
359 Westerdijk Institute under number NCCB100732
360 (<http://www.westerdijkinstituut.nl/Collections/DefaultInfo.aspx?Page=Bacteria>), the National
361 Collection of Type Cultures of Public Health England under number NCTC14359
362 (<https://www.phe-culturecollections.org.uk/collections/nctc.aspx>) or through the
363 corresponding author.

364

365 Appendixes

366 **The COMBAT consortium (in alphabetical order)**

367 Maris S. Arcilla¹, Martin C.J. Bootsma², Perry J. van Genderen³, Abraham Goorhuis⁴, Martin
368 Grobusch⁴, Jarne M. van Hattem⁵, Menno D. de Jong⁵, Damian C. Melles¹, Nicky Molhoek⁶,
369 Astrid M.L. Oude Lashof⁷, John Penders⁷, Constance Schultsz^{5, 8}, Ellen E. Stobberingh⁵, Henri A.
370 Verbrugh¹

371 ¹Erasmus University Medical Centre, Department of Medical Microbiology and Infectious
372 Diseases, Rotterdam,

373 ²Utrecht University, Department, Faculty of Science, Utrecht,

374 ³Havenziekenhuis - Institute for Tropical Diseases, Department of Internal Medicine,
375 Rotterdam,

376 ⁴Academic Medical Centre, Center of Tropical Medicine and Travel Medicine, Amsterdam,

377 ⁵Academic Medical Centre, Department of Medical Microbiology, Amsterdam,

378 ⁶Havenziekenhuis - Institute for Tropical Diseases, Travel Clinic, Rotterdam,

379 ⁷Maastricht University Medical Centre, School for Nutrition, Toxicology and Metabolism and
380 School for Public Health and Primary Care, Department of Medical Microbiology, Department of
381 Medical Microbiology, Maastricht,

382 ⁸Academic Medical Centre, Department of Global Health-Amsterdam Institute for Global Health
383 and the Development, Amsterdam,

384 All in the Netherlands

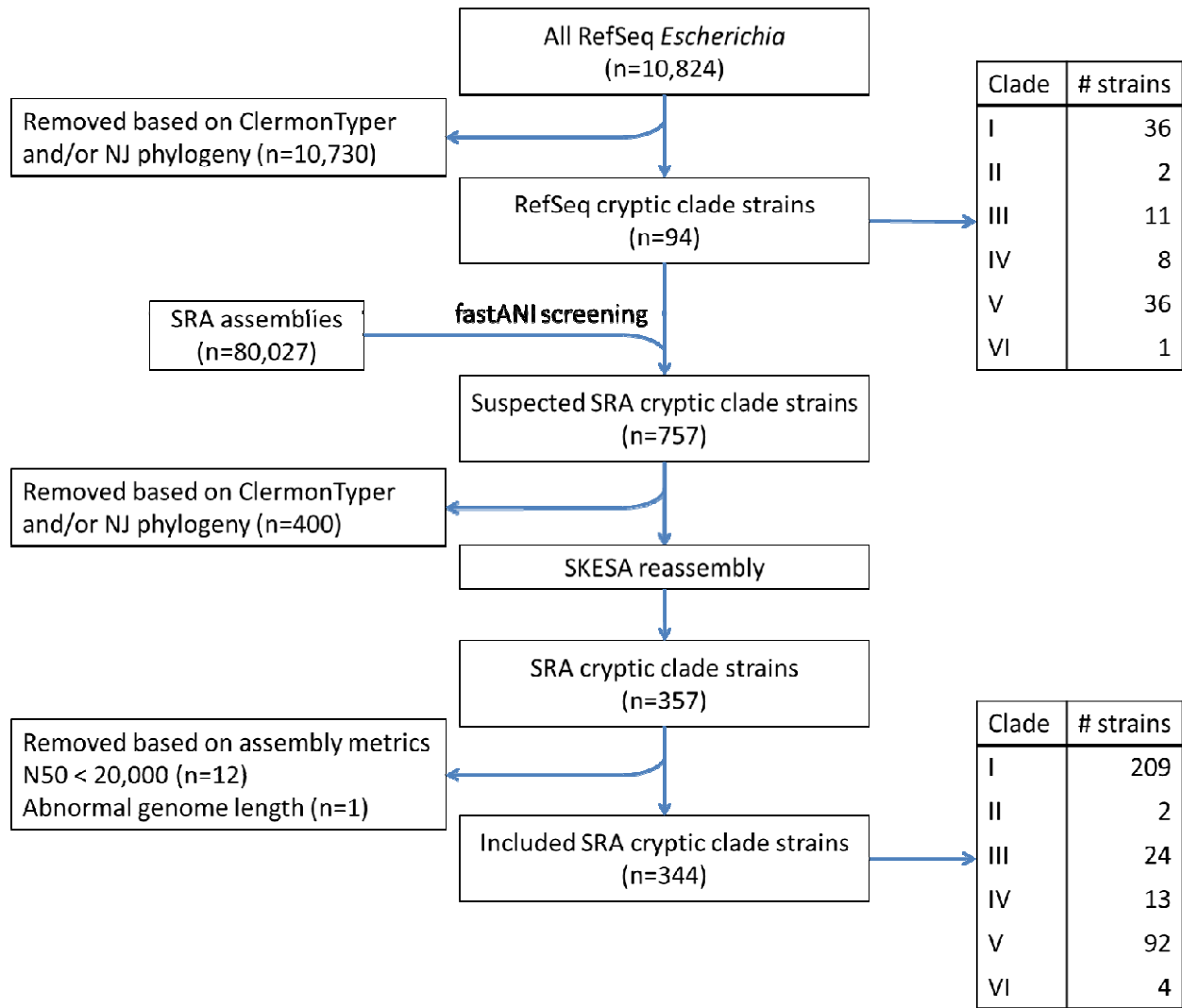
385

386 References

- 387 1. Farmer JJ, Fanning GR, Davis BR, et al. *Escherichia fergusonii* and *Enterobacter taylorae*, two new
388 species of Enterobacteriaceae isolated from clinical specimens. *J Clin Microbiol.* 1985;21(1):77-81.
- 389 2. Huys G, Cnockaert M, Janda JM, Swings J. *Escherichia albertii* sp. nov., a diarrhoeagenic species
390 isolated from stool specimens of Bangladeshi children. *Int J Syst Evol Microbiol.* 2003;53(3):807-
391 810.
- 392 3. Liu S, Jin D, Lan R, et al. *Escherichia marmotae* sp. nov., isolated from faeces of *Marmota*
393 *himalayana*. *Int J Syst Evol Microbiol.* 2015;65(7):2130-2134.
- 394 4. Russo TA, Johnson JR. Medical and economic impact of extraintestinal infections due to
395 *Escherichia coli*: Focus on an increasingly important endemic problem. *Microbes Infect.*
396 2003;5(5):449-456. doi:10.1016/S1286-4579(03)00049-2
- 397 5. Savini V, Catavittello C, Talia M, et al. Multidrug-Resistant *Escherichia fergusonii*: a Case of Acute
398 Cystitis. *J Clin Microbiol.* 2008;46(4):1551-1552. doi:10.1128/JCM.01210-07
- 399 6. Ooka T, Seto K, Kawano K, et al. Clinical Significance of *Escherichia albertii*. *Emerg*
400 *Infect Dis J.* 2012;18(3):488. doi:10.3201/eid1803.111401
- 401 7. Gangiredla J, Mammel MK, Barnaba TJ, et al. Draft Genome Sequences of *Escherichia albertii*,
402 *Escherichia fergusonii*, and Strains Belonging to Six Cryptic Lineages of *Escherichia* spp. *Microbiol*
403 *Resour Announc.* 2018;6(18). doi:10.1128/genomeA.00271-18
- 404 8. Walk ST, Alm EW, Gordon DM, et al. Cryptic lineages of the genus *Escherichia*. *Appl Environ*
405 *Microbiol.* 2009;75(20):6534-6544. doi:10.1128/AEM.01262-09
- 406 9. Beghain J, Bridier-Nahmias A, Le Nagard H, Denamur E, Clermont O. ClermonTyping: an easy-to-
407 use and accurate in silico method for *Escherichia* genus strain phylotyping. *Microb Genomics.*
408 2018:1-8. doi:10.1099/mgen.0.000192
- 409 10. Ondov BD, Treangen TJ, Melsted P, et al. Mash: Fast genome and metagenome distance
410 estimation using MinHash. *Genome Biol.* 2016;17(1):1-14. doi:10.1186/s13059-016-0997-x
- 411 11. Simonsen M, Mailund T, Pedersen CNS. Rapid Neighbour-Joining. In: Crandall KA, Lagergren J, eds.
412 *Algorithms in Bioinformatics*. Berlin, Heidelberg: Springer Berlin Heidelberg; 2008:113-122.
- 413 12. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUASt: Quality assessment tool for genome
414 assemblies. *Bioinformatics.* 2013;29(8):1072-1075. doi:10.1093/bioinformatics/btt086
- 415 13. Ochman H, Selander RK. Standard reference strains of *Escherichia coli* from natural populations. *J*
416 *Bacteriol.* 1984;157(2):690-693.
- 417 14. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective stochastic
418 algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol.* 2015;32(1):268-274.
419 doi:10.1093/molbev/msu300

- 420 15. Gardner SN, Slezak T, Hall BG. kSNP3.0: SNP detection and phylogenetic analysis of genomes
421 without genome alignment or reference genome. *Bioinformatics*. 2015;31(17):2877-2878.
422 doi:10.1093/bioinformatics/btv271
- 423 16. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact
424 alignments. *Genome Biol*. 2014;15:R46. doi:10.1186/gb-2014-15-3-r46
- 425 17. Chun J, Oren A, Ventosa A, et al. Proposed minimal standards for the use of genome data for the
426 taxonomy of prokaryotes. *Int J Syst Evol Microbiol*. 2018;68(1):461-466.
427 doi:10.1099/ijsem.0.002516
- 428 18. Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM. DNA–DNA
429 hybridization values and their relationship to whole-genome sequence similarities. *Int J Syst Evol
430 Microbiol*. 2007;57(1):81-91.
- 431 19. Pettengill EA, Pettengill JB, Binet R. Phylogenetic Analyses of Shigella and Enteroinvasive
432 Escherichia coli for the Identification of Molecular Epidemiological Markers: Whole-Genome
433 Comparative Analysis Does Not Support Distinct Genera Designation. *Front Microbiol*.
434 2016;6:1573. doi:10.3389/fmicb.2015.01573
- 435 20. Arcilla MS, van Hattem JM, Haverkate MR, et al. Import and spread of extended-spectrum beta-
436 lactamase-producing Enterobacteriaceae by international travellers (COMBAT study): a
437 prospective, multicentre cohort study. *Lancet Infect Dis*. 2017;17(1):78-85. doi:10.1016/S1473-
438 3099(16)30319-X
- 439 21. McArthur AG, Waglechner N, Nizam F, et al. The comprehensive antibiotic resistance database.
440 *Antimicrob Agents Chemother*. 2013;57(7):3348-3357. doi:10.1128/AAC.00419-13
- 441 22. Chen L, Yang J, Yu J, et al. VFDB: A reference database for bacterial virulence factors. *Nucleic Acids
442 Res*. 2005;33(DATABASE ISS.):325-328. doi:10.1093/nar/gki008
- 443 23. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI analysis of
444 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun*. 2018;9(1):5114.
445 doi:10.1038/s41467-018-07641-9
- 446 24. Rodriguez-R LM, Konstantinidis KT. The enveomics collection: a toolbox for specialized analyses of
447 microbial genomes and metagenomes. *PeerJ Prepr*. 2016;4:e1900v1.
448 doi:10.7287/peerj.preprints.1900v1
- 449 25. Yoon S-H, Ha S-M, Lim J, Kwon S, Chun J. A large-scale evaluation of algorithms to calculate
450 average nucleotide identity. *Antonie Van Leeuwenhoek*. 2017;110(10):1281-1286.
451 doi:10.1007/s10482-017-0844-4
- 452 26. Walk S. The “Cryptic” Escherichia. *EcoSal Plus*. 2015.
453 <https://www.asmscience.org/content/journal/ecosalplus/10.1128/ecosalplus.ESP-0002-2015>.
- 454 27. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. Selection of
455 representative genomes for 24,706 bacterial and archaeal species clusters provide a complete
456 genome-based taxonomy. *bioRxiv*. 2019. doi:10.1101/771964

- 457 28. Alikhan N-F, Bachmann NL, Zakour NL Ben, et al. Multiple evolutionary trajectories for non-O157
458 Shiga toxigenic *Escherichia coli*. *bioRxiv*. 2019. doi:10.1101/549998
- 459 29. Lees JA, Galardini M, Bentley SD, Weiser JN, Corander J. pyseer: a comprehensive tool for microbial
460 pangenome-wide association studies. *Bioinformatics*. 2018;34(24):4310-4312.
461 doi:10.1093/bioinformatics/bty539
- 462 30. Souvorov A, Agarwala R, Lipman DJ. SKESA: strategic k-mer extension for scrupulous assemblies.
463 *Genome Biol*. 2018;19(1):153. doi:10.1186/s13059-018-1540-z
- 464 31. Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS. UFBoot2: Improving the Ultrafast
465 Bootstrap Approximation. *Mol Biol Evol*. 2017;35(2):518-522. doi:10.1093/molbev/msx281
- 466 32. Csardi G, Nepusz T. The igraph software package for complex network research. *InterJournal*.
467 2006;Complex Sy:1695.
- 468 33. Ooms J. The jsonlite Package: A Practical and Consistent Mapping Between JSON Data and R
469 Objects. *ArXiv E-Prints*. 2014:arXiv:1403.2805.
- 470 34. Shannon P, Markiel A, Ozier O, et al. Cytoscape: a software environment for integrated models of
471 biomolecular interaction networks. *Genome Res*. 2003;13(11):2498-2504. doi:10.1101/gr.1239303
- 472 35. Zhou Y, Chen Y, Chen S, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*.
473 2018;34(17):i884-i890. doi:10.1093/bioinformatics/bty560
- 474 36. Meier-Kolthoff JP, Auch AF, Klenk H-P, Göker M. Genome sequence-based species delimitation
475 with confidence intervals and improved distance functions. *BMC Bioinformatics*. 2013;14(1):60.
476 doi:10.1186/1471-2105-14-60
- 477 37. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new developments. *Nucleic
478 Acids Res*. 2019;47(W1):W256-W259. doi:10.1093/nar/gkz239
- 479
- 480

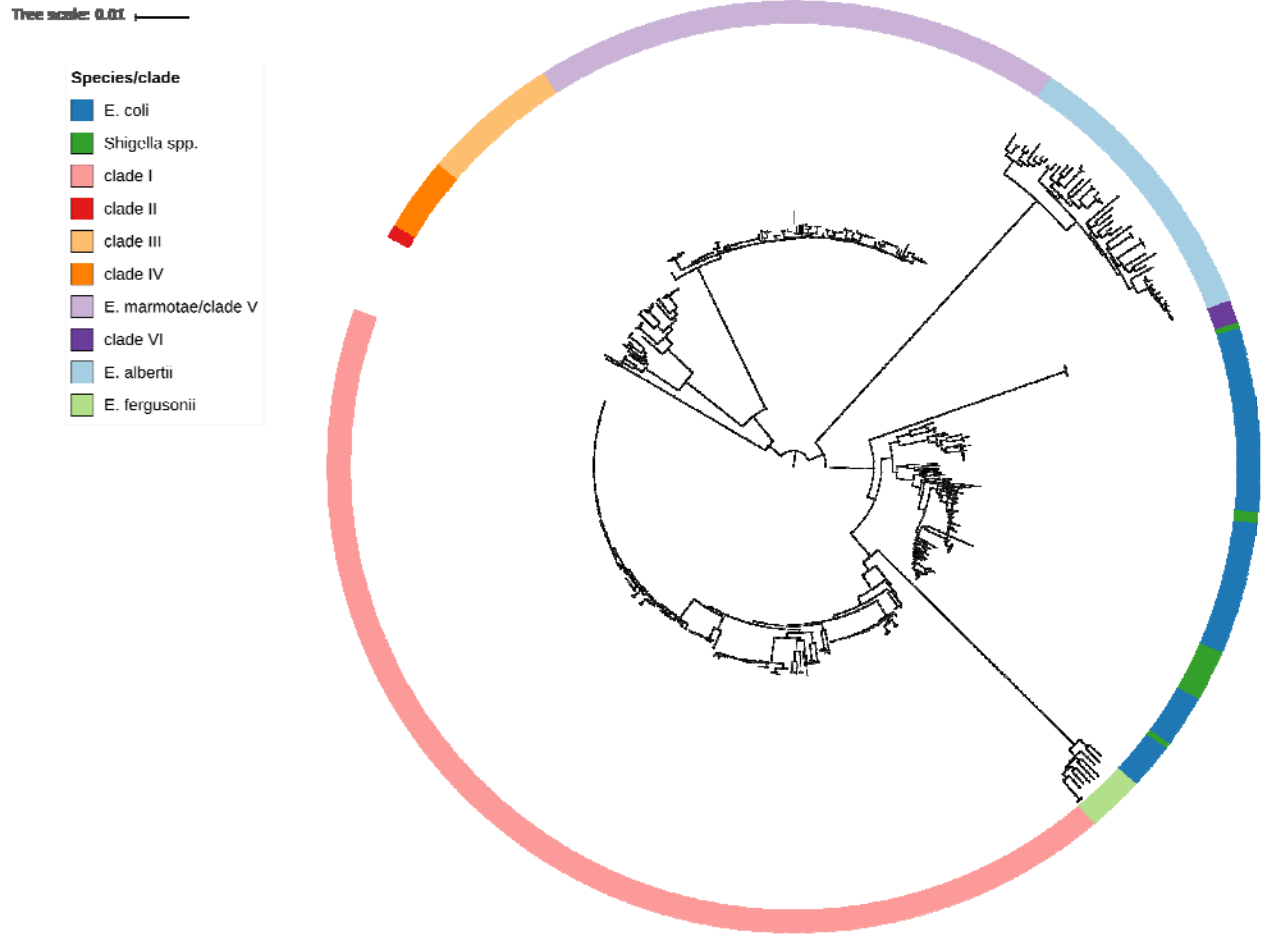


481

482 **Figure 1.** Flowchart of inclusion process of cryptic clade strains. Only the well-assembled SRA

483 cryptic clade strains (lower table on right hand side) are used in subsequent analyses.

484



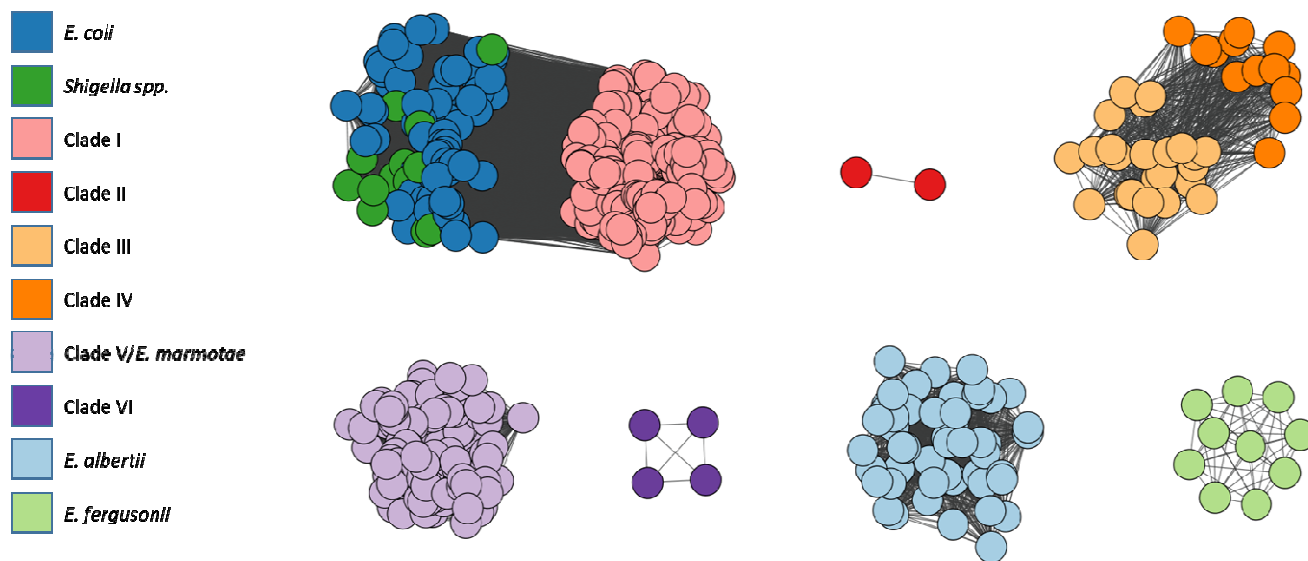
485

486 **Figure 2.** Midpoint-rooted maximum-likelihood phylogeny of 489 *Escherichia* and *Shigella*

487 strains based on a core SNP alignment of 3445 positions. Nodes with bootstrap values under 95

488 were collapsed. Visualized in iTOL³⁷.

489



490

491 **Fig 3.** Network of 489 included *Escherichia* and *Shigella* strains, based on fastANI analysis. Edges

492 represent ANI values over 95% (boundary for species delineation) between two strains.

493