1   Ecotype-specific genomic features within the *Escherichia* cryptic clade IV

2   Martín Saraceno[a,b], Nicolás Frankel[b,c], Martín Graziano[a,b]∗

3   martinsa@ege.fcen.uba.ar, nfrankel@ege.fcen.uba.ar, marting@ege.fcen.uba.ar

4   [a]CONICET - Universidad de Buenos Aires. Instituto de Ecología, Genética y Evolución
5   de Buenos Aires (IEGEBA). 2160 Intendente Güiraldes St., Buenos Aires-C1428EGA,
6   Argentina. [b]Universidad de Buenos Aires. Facultad de Ciencias Exactas y Naturales.
7   Departamento de Ecología, Genética y Evolución. 2160 Intendente Güiraldes St.,
8   Buenos Aires-C1428EGA, Argentina. [c]CONICET - Universidad de Buenos Aires.
9   Instituto de Fisiología, Biología Molecular y Neurociencias (IFIBYNE). 2160
10  Intendente Güiraldes St., Buenos Aires-C1428EGA, Argentina.

11  *Corresponding author: marting@ege.fcen.uba.ar

12  **Abstract**

13  *Escherichia* cryptic clades represent a relatively unexplored taxonomic cluster believed
14  to exhibit characteristics associated with a free-living lifestyle, which is known as the
15  environmental hypothesis. This hypothesis suggests that certain *Escherichia* strains
16  harbour traits that favour their environmental persistence, thus expanding the ecological
17  commensal niche of the genus. While surveying *Escherichia* diversity in an urban South
18  American stream we isolated the first environmental cryptic clade IV strain in South
19  America (339_SF). Here we report the genomic characterization of 339_SF strain in the
20  context of existing genomic information for cryptic clade IV. A comparative analysis of
21  genomes within the same clade stemming from diverse ecological sources and
22  geographical locations reveals close phylogenetic proximity between our isolate and
23  strains of environmental origin.  In the genomes of cryptic clade IV strains that were
24  isolated from environmental niches we observed enrichment of functional genes related
25  to responses to adverse environmental conditions and a low number of genes with
26  clinical relevance among. Our findings highlight substantial intra-group genomic
27  structuring linked to ecological origin and shed light on the genomic mechanisms
28  underlying the naturalization phenomena within the *Escherichia* genus.

29  **Keywords**

30  phylogenomics;   environmental   hypothesis;   microbial   free-living   ecology;
31  environmental genomics.

32  **Introduction**

33  The *Escherichia* genus includes widely known species such as *E. coli*, *E. fergusonii* and
34  *E. albertii,* and five monophyletic groups named "cryptic" clades I to V (Walk, 2015).
35  The name of the latter is based on the inability to distinguish them from representative
36  isolates of *E. coli* through typical biochemical diagnostic reactions. However, relatively
37  recent phylogenomic analyses have shown that these lineages are divergent from stem
38  members of the genus (Walk *et al.*, 2009; Luo *et al.*, 2011). Previous studies associated
39  some cryptic clades with a free-living lifestyle, a phenomenon known as the
40  *environmental hypothesis*. The initial evidence for this proposition was based on the
41  overrepresentation of these clades in environmental samples (Walk *et al.*, 2007, 2009).

42    Subsequent phenotypic, genomic and transcriptomic analyses supported the link as well
43    (Walk, 2015; Di Sante *et al.*, 2018), but the evidence behind this hypothesis is still limited.

44    Cryptic clade I is closely related to *E. coli* and genomic studies indicate that both groups
45    carry common virulence factors, which suggests that clade I strains could also be as
46    pathogenic as some *E. coli* strains (Steinsland *et al.*, 2010; Walk, 2015). Bacterial virulence
47    factors encode genes that facilitate infection, such as toxins or proteins needed for
48    bacterial adherence (Holden *et al.*, 2009; Stecher and Hardt, 2011; Acosta-Dibarrat *et al.*,
49    2021). The presence of virulence and antibiotic resistance factors in bacterial genomes is
50    often linked to specific selection pressures associated with hosts (Jernberg *et al.*, 2010;
51    Becattini *et al.*, 2016). On the other hand, field, experimental and genomic evidence
52    support the hypothesis that the remaining clades (II to V) survive outside animal hosts
53    (Di Sante *et al.*, 2018). This assumption is based not only on the fact that they have been
54    generally isolated from soil or surface water (they have not been linked to cases of
55    human or animal infection), but also on the lack of virulence factors for intestinal and
56    extra-intestinal infection (Ingle *et al.*, 2011; Vignaroli *et al.*, 2015). Moreover, a
57    comparative genomic study suggested that genomes of cryptic clades associated with a
58    free-living style are enriched in genes that confer greater fitness in the environment, e.g.
59    contributing to novel metabolic pathways for the exploitation of alternative energy
60    sources (Luo *et al.*, 2011). A caveat of the analyses mentioned above is that they were
61    performed with only a few genomes. Thus, a more comprehensive analysis that includes
62    all currently available genomes is in need for a better understanding of the genomic
63    signatures related to a free-living style.

64    Addressing the ecotype-genetic relationship within the genus *Escherichia* is relevant
65    due to the increasing number of studies that have reported the ability of certain *E. coli*
66    strains to persist in secondary habitats (Ishii *et al.*, 2006; Lee *et al.*, 2006; Mackowiak *et al.*,
67    2018). Thus, the review of the traditional niche assigned to *E. coli* is framed by a larger
68    discussion on the ecological history of the genus itself, which calls for an increase in
69    genomic and phenotypic characterizations. It has been hypothesized that the
70    phenomenon of environmental persistence could be the result of a complex balance
71    between the differential fitness of some *E. coli* strains and the occurrence of permissive
72    ecological conditions, i.e., optimal physicochemical characteristics, nutrient availability,
73    low competition, among others (Surbeck *et al.*, 2010; Jang *et al.*, 2017). Furthermore,
74    secondary habitats such as urban streams, where point and diffuse sources of bacterial
75    input usually coexist with the native microbiome and temperatures are usually benign
76    and nutrient availability high, are hot spots for genetic exchange and the emergence of
77    new strains, which can be pathogenic (McLellan *et al.*, 2015). In this line, previous
78    research has detected the imprint of isolation sources both in genomic traits of
79    epidemiological importance, as well as in phenotypes linked to long-term persistence in
80    *E. coli* (Berthe *et al.*, 2013). In a recent study, also on *E. coli*, associations among genetic
81    backgrounds and specific habitats were uncovered and horizontal gene flow was found
82    to be an important mechanism driving the reinforcement of gene structuring (Touchon *et*
83    *al.*, 2020). Due to the sanitary and ecological interest of these phenomena, we consider
84    that it is relevant to explore the genomic features associated with a free-living style
85    within the genus *Escherichia*.

86 In a previous work, we reported the isolation of a member of cryptic clade IV from a
87 South American urban stream (Saraceno *et al.*, 2020). Here we report the genomic
88 characterization of this strain, while performing a comparative analysis of genomes
89 within the clade IV from diverse ecological and geographic sources. We present
90 evidence of intra-clade IV functional genes structuring linked to the ecotype of origin.
91 We discuss our results within the framework of the environmental hypothesis and the
92 occurrence of niche-specific selective pressures.

93 **Methods and Materials**

94 *Environmental cryptic clade IV isolation and phylogenetic assignment*

95 The isolate 339_SF, belonging to the cryptic clade IV, was isolated, cryopreserved and
96 phylogenetically characterized as previously mentioned in Saraceno *et al.* (2020). Briefly,
97 environmental isolates were obtained from the water column of San Francisco urban
98 stream (Buenos Aires, Argentina) by culturing methods: a first-round employing
99 Chromocult® Coliform Agar selective medium (MilliporeSigma) and, in a second step,
100 blue- to purple-coloured selected colonies were streaked at least two times onto Levine
101 E.M.B. (Eosin Methylene Blue) agar plates to assure isolates purity. The phylogenetic
102 assignment was carried out through a series of multiplex PCR procedures. A first round,
103 through the amplification of the *araA*, *chuA*, *yjaA* and *TspE4.C2* genes, which
104 determined that isolate 339_SF belonged to an *Escherichia* cryptic clade, and a
105 subsequent second round of a double PCR method, based on the amplification of *aes*
106 and *chuA* genes, to finally determined its cryptic lineage membership (Clermont *et al.*,
107 2011, 2013). All PCR procedures were carried out using a loaded loop after ringing on
108 each colony as a template, in a 20 µL final volume, and 10 µL of 2X GoTaq® Green
109 Master Mix (Promega), following respective literature indications.

110 *DNA extraction and sequencing*

111 Before DNA extraction, 339_SF isolate was cultured overnight in Luria Bertani broth at
112 37°C and 1,5 ml of this culture was pelleted by ultracentrifugation (13000 rpm for 5
113 min). Obtained pellet was incubated with saline EDTA buffer ($C_{10}H_{16}N_2O_8$ 0.01 M and
114 NaCl 0.15 M; pH = 8.0) and proteinase K. DNA purification was held through serial
115 washes with chloroform:isoamyl alcohol (24:1), precipitated with isopropyl alcohol and
116 re-suspended in low-TE (Tris 1 mM and EDTA 0.1 mM; pH = 7.0). Genomic DNA
117 integrity was assessed through agarose gel electrophoresis (1% agarose). DNA mass
118 was estimated by the inclusion of a mass ladder (MassRulerTM DNA Ladder Mix) and
119 its concentration by a Qubit assay (dsDNA Quantitation, broad range). Whole genome
120 sequencing was performed in Novaseq platform (Illumina, 2x150 paired-end), obtaining
121 a Q30 of 91.99% and an estimated sequencing coverage of 890X.

122 *Genome assembly and gene annotation*

123 Reads quality was assessed using FASTQC and Kraken2 (de Sena Brandine and Smith,
124 2019; Wood *et al.*, 2019). Genome assembly was made using Unicycler (Wick *et al.*, 2017).
125 The quality of the assembly was assessed by Quast (Gurevich *et al.*, 2013) and its
126 phylotyping was corroborated using Clermontyping *in silico* method (Beghain *et al.*,
127 2018). N50 was of 536699 bp and L50 of 3 contigs, revealing a successful assembly
128 procedure. Gene annotation was done using Prokka (minimum contig size 300 bp,
129 genus-specific BLAST databaseforg. *Escherichia*) (Seemann, 2014). Search for genome

130 features of clinical and epidemiological interest was done using ABRicate tool and
131 VFDB database for virulence factors, NCBI and Resfinder databases for antimicrobial
132 resistance (AMR) genes and Plasmidfinder database for plasmids replicons (Carattoli *et*
133 *al.*, 2014; Seemann, 2020; Florensa *et al.*, 2022; Liu *et al.*, 2022), considering only those with
134 both coverage and identity above 95%.

135

136 *Phylogenetics*

137 Genome assemblies were obtained from public and curated databases such as
138 EnteroBase (Zhou *et al.*, 2020), belonging to the cryptic clade IV (44 genomes) and the
139 cryptic clades I (2 genomes), II (1 genome), III (2 genomes) and V (2 genomes). Also,
140 genome assemblies from *E. albertii*, *E. fergusonni*, *E. coli* K-12 and *E. coli* main
141 pathogenic groups were acquired: EPEC (enteropathogenic *E. coli*), ETEC
142 (enterotoxigenic *E. coli*), EIEC (enteroinvasive *E. coli*), STEC (Shiga toxin-producing
143 *E. coli*), EAEC (enteroaggregative *E. coli*), EHEC (enterohemorrhagic *E. coli*), UPEC
144 (uropathogenic *E. coli*) y NMEC (neonatal meningitis-causing *E. coli*). Among clade IV
145 genomes, 10 belong to the environmental isolation ecotype, 14 to human hosts, and 20
146 to other hosts (wild and domestic animals); all of them from different geographic
147 locations around the world. Accession information and metadata from the 62 genomes
148 are listed in **Table S1**. All the assemblies were analyzed following the same paths of
149 annotation and search for genomic features of clinical interest mentioned above. From
150 the annotation files obtained, a pangenome analysis was conducted through Roary
151 pipeline (Page *et al.*, 2015), identifying the core and accessory genes among genomes.
152 Core genes were employed as input for the construction of the phylogenetic
153 relationships. The maximum likelihood method with a bootstrapping of 1000 was
154 employed for the inference of the tree using IQ-TREE (Nguyen *et al.*, 2015). Lastly, the
155 resulting phylogenetic tree was rooted in the midpoint and the graphics were edited
156 using FigTree and Inkscape, respectively (Harrington, 2004; Rambaut, 2018).

157 *Pangenome of cryptic clade IV and genetic enrichment analysis according to ecotype*

158 The set of annotated genomes belonging to the cryptic clade IV (44) was used as input
159 to design a presence/absence matrix of functional genes - 3524 entries-, as characterized
160 by UniProtKB database (The Consortium Uniprot, 2023). Genes detected in only one
161 genome or across the entire genome set were discarded, and a hierarchical clustering
162 analysis was held with the resulting matrix. Further, the set was subdivided into two
163 subsets according to its ecotype host/environmental origin. The number of times each
164 gene was detected in each group was counted (repeated copies of the same gene per
165 genome were considered only one time). Consequently, the detection frequencies of
166 each gene were calculated inside each subset. To consider a gene as enriched within one
167 of the two ecotypes, a mixed criterion was used based on the resulting frequencies: a
168 minimum threshold of 80% detection in one of the subsets was required and, jointly, a
169 detection frequency of less than 80% in the complementary subset. After the
170 identification of the genes considered enriched in each ecotype, a review of the
171 literature was carried out to assign them a functional and ecological context.

172 **Results**

173 *Genomic features of the environmental cryptic clade IV isolate*

174   Strain 339_SF was isolated from an urban stream in South America and characterized as
175   a member of the cryptic clade IV (Saraceno *et al.*, 2020). To our knowledge, there is no
176   previous report of an environmental isolate belonging to *Escherichia* cryptic clades in
177   the region. We sequenced the whole genome of strain 339_SF and produced a total of
178   29 contigs with 4283 predicted coding sequences. The total length of the assembly
179   (4637586 bp) and the percentage of GC content (50.86%) are consistent with the
180   expectations for an *Escherichia* genome (Hildebrand *et al.*, 2010). No genes associated
181   with antibiotic resistance or plasmid replicons were detected, while a total of 16
182   sequences matched virulence factors (VFs). Among the VFs detected, none are strictly
183   associated with a defined *E. coli* pathotype. For example, we found factor *sat*, which is
184   often associated with UPEC and EAEC profiles, but which is distributed between
185   commensal strains as well (Toloza *et al.*, 2015). In addition, we detected the numerous
186   variants of the type 1 fimbria precursor gene (*fimB*, *fimC*, *fimD, fimG* and *fimI*), which
187   are common VFs in strains associated with UPEC infections (Müller *et al.*, 2009). We
188   also identified factors associated with siderophore function (*chuX,fepA, fepB, fepD*and
189   *entS*), which are widespread among various *E. coli* strains and with presumed negative
190   effects on hosts (Ozenberger *et al.*, 1987; Bleuel *et al.*, 2005; Suits *et al.*, 2009) and factors
191   that encode components of the general gram-negative secretion pathway (*gspG*, *gspH*,
192   *gspI* and *gspL*) (Py *et al.*, 2001). *KpsD*, a key factor for the synthesis of the capsule in *E.*
193   *coli,* which is the structure that confers resistance during extraintestinal infections
194   (Russo *et al.*, 1998; Duan *et al.*, 2020), was also detected.

195   *Genes of clinical interest in Clade IV genomes*

196   We compared available genomes of cryptic clade IV with that of 339_SF. Among the 44
197   genomes, the minimum number of VFs detected was 10 and the maximum was 23, with
198   an average of 14.8. Of the 43 VFs found across the set, 12 were shared by 84.4% of the
199   clade IV genomes. These 12 factors are present in the genome of 339_SF (see **Table S2**
200   for the genomic traits of clinical and epidemiological relevance). The genomes carrying
201   more VFs were isolated from humans from different continents (ESC_BA8399AA,
202   ESC_EA6501AA and ESC_LA2398AA, with 23 putative genes each, respectively from
203   Africa, Europe and South America). We did not detect AMR genes in the 339_SF
204   genome, as in most of the cryptic clade IV genomes (39 of 45 genomes have no
205   resistance determinants). However, we detected 13 factors in the genome of a poultry
206   isolate (ESC_FA7484AA). Similar results are observed when looking at plasmid
207   replicons: clade IV genomes carry one plasmid replicon on average, and almost half of
208   them do not carry replicons. Thus, 339_SF possesses a similar profile of factors of
209   clinical interest when compared to other cryptic clade IV genomes from various
210   geographic or ecotype origins.

211   *Phylogenomics of Cryptic Clade IV*

212   We performed a phylogenetic analysis based on *core* genes with the available genomes
213   of clade IV, genomes of reference strains from the rest of the cryptic clades and other
214   main taxa of the genus *Escherichia.* After annotating genes across all assemblies, 501
215   were found in the entire set (**Table 1**). Clade IV strains formed a monophyletic group in
216   the tree (**Figure 1**). Genomes belonging to clades II, III and V appear close to clade IV.
217   Clade I genomes fall close to *E. coli* strains (K-12 and those considered pathogenic) and
218   *E. fergusonii*. 339_SF groups with a subset of clade IV genomes that include 6 other

219 isolates of environmental origin. Furthermore, this cluster of genomes tends to lack
220 antibiotic resistance genes and plasmid replicons, while carrying a similar number of
221 VFs - between 14 and 16- (**Table S2**).

222 Table 1) Abundance of gene types according to their sharing percentage between the strains.

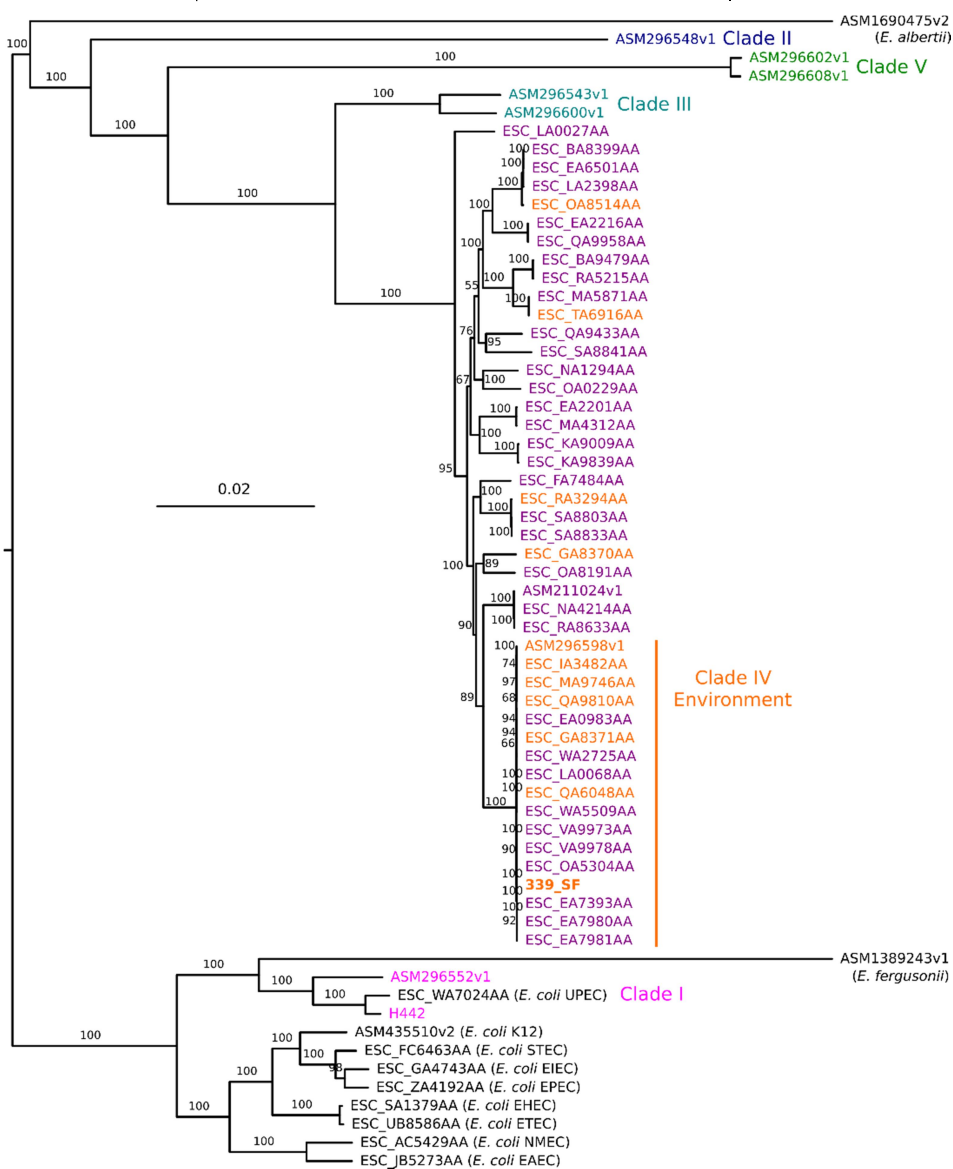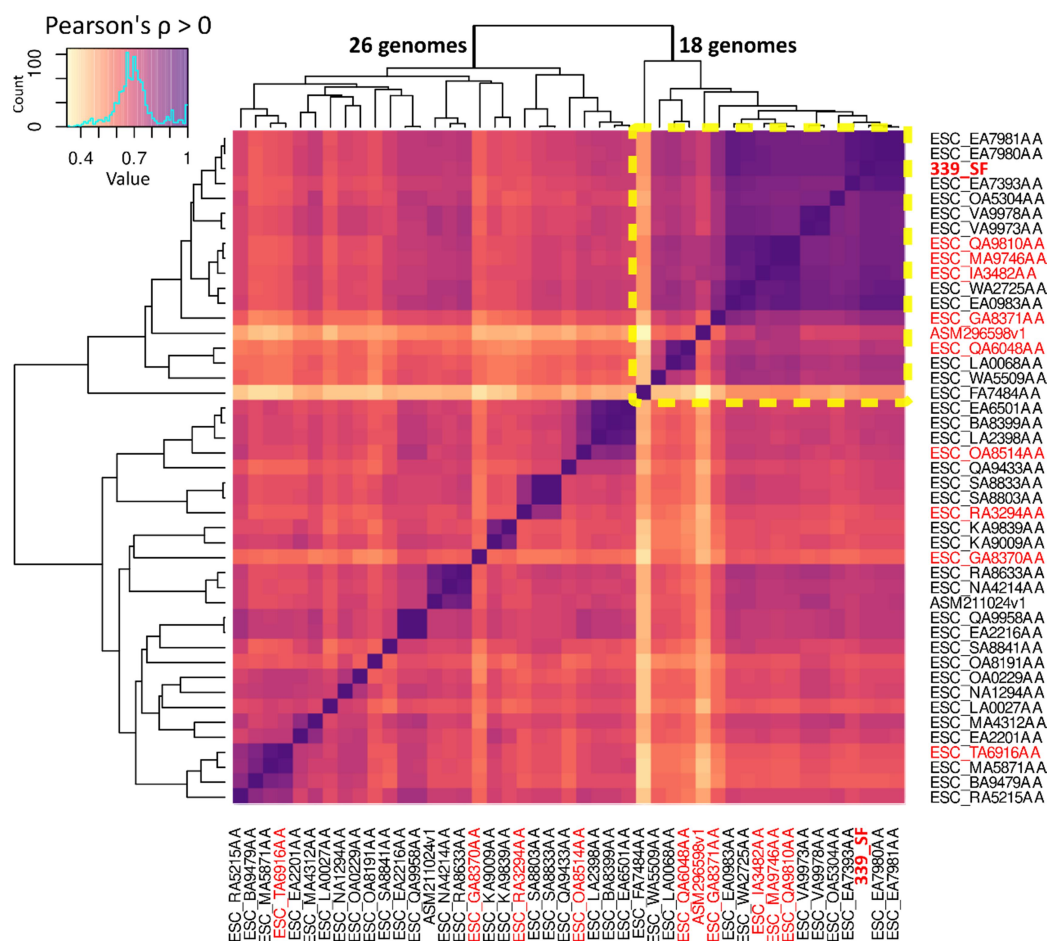| Types of genes | Percentage sharing between strains | Number of genes |
|---|---|---|
| *core* | 99% <= genomes <= 100% | **501** |
| *soft core* | 95% <= genomes <= 99% | 1604 |
| *shell* | 15% <= genomes <= 95% | 2493 |
| *cloud* | 0% <= genomes <= 15% | 18485 |
| Total | 0% <= genomes <= 100% | **23083** |



223

224 **Figure 1**: Phylogenetic relationships amongst selected *Escherichia* genus strains. Genomes from cryptic
225 clades IV (ecotypes are discriminated by colour, host -orange- or environmental -violet-), I, II, III and V,

226  E. coli K12 and reference strains from its main pathotypes, as well as *E. fergusonii* and *E. albertii* were
227  included. 339_SF is labeled in bold orange and the cluster containing relatively more environmental
228  ecotype genomes is indicated. Full sequences from the 501 core genes were aligned presenting a base
229  overlap region of 0.4 Mb and phylogeny was conducted as detailed in Materials and Methods. Bootstrap
230  support values are indicated on each branch.

231  *Association between genomic features and ecotype within cryptic clade IV*

232  We grouped genomes following a functional genetic content perspective to further
233  investigate the relationship between genomic features and ecotype within cryptic clade
234  IV. Hence, a gene presence/absence matrix was created using a pangenome (the totality
235  of annotated genes with experimentally validated or inferred functionality across all the
236  set of clade IV genomes). After purging the matrix of genes detected in a single genome
237  and present in all genomes, we retained 870 entries. A hierarchical clustering based on
238  Pearson's positive correlations was performed with these data  (**Figure 2**). The tree
239  shows that CC IV genomes are grouped into two main clusters of 18 and 26 genomes.
240  Globally, this organization is similar to that observed in the phylogeny based on core
241  genes (**Figure 1**). 339_SF is included in the group that contains the highest proportion
242  of isolates of the environmental ecotype (7 of 18 genomes, versus 4 of 26). We also
243  observed that 13 of the genomes within the cluster of 18 have a very strong positive
244  correlation (upper right corner, **Figure 2**). 339_SF and 4 other genomes from
245  environmental origin are part of these highly correlated genomes.



246

247 **Figure 2**: **Hierarchical clustering of cryptic clade IV genomes based on functional genes**
248 **presence/absence**. The colours of the heat map indicate the Pearson correlation coefficient between 0 to
249 1 among the genomes. The colour of the label indicates the ecotype: red for environmental origin
250 genomes and black for host origin. The cluster containing the 339_SF, labelled in bold red, is boxed in a
251 yellow dotted line.

252 In parallel, we looked for genes enriched in clade IV genomes in relation to their
253 ecotype (environmental or host). We found 24 over-represented genes in environmental
254 genomes and only 3 over-represented genes in host genomes (**Table S3**). Of the 24
255 genes enriched in environmental ecotypes, 12 were previously described for *E. coli*
256 reference strains, while the remaining 12 were identified in reference genomes
257 belonging to other bacterial species. These genes may be linked to different key
258 biological functions related to a free-living style. A group of factors are directly
259 associated with the response to stress conditions and DNA damage: *UmuC* and *UmuD*
260 help repair DNA damage and participate in the SOS response, *yrecE* and *ybcO* are
261 nucleases that intervene in DNA repair, *dnaK* has chaperone activity during stress
262 events and *kilR* encodes an inhibitor of cell division in response to antibiotics. The gene
263 *hicA*, which encodes for a component of a type II toxin-antitoxin system in *E. coli*, has
264 also been associated with the stress response. In addition, genes related to viral
265 infections were also enriched in the environmental set: *cas6f* and *csy3* belong to the
266 CRISPR defence system, *hpaIIM* encodes a restriction enzyme described in
267 *Haemophilus parainfluenzae*, and *intQ*, is related to the integration of phages into
268 bacterial genomes. Other enriched genes are the *bepC* genes, which confer antibiotic
269 resistance to substances of hydrophobic nature, and *icsA*, which encodes a protein
270 essential for bacterial adhesion and virulence (described in *Shigella flexneri*). Also, a
271 group of genes associated with the environmental ecotype is linked to structural and
272 transmembrane transport functions (*yidK*, *csbX*, *ompC*, *yidI*, *yknY*, *fliC* and *gpFI*).
273 Finally, genes linked to energetic metabolism were also identified (*ahr*, which is
274 involved in lipid metabolism, and *hypBA1*, related to carbohydrate metabolism).

275 On the other hand, over-represented genes among the host-origin genomes were *rrrD*,
276 which encodes a factor with hydrolytic activity of bacterial walls, and *tolA*, which
277 encodes a Tol-Pal system factor carrying key functions in cell division and outer
278 membrane integrity. Both genes were previously described in *E. coli* reference strains. A
279 third gene associated with host ecotypes is *prfA*, which was described in *Mycobacterium*
280 *tuberculosis* as involved in protein biosynthesis.

281 **Discussion**

282 Phylogenetic analyses combined with comparative genomics approaches can help
283 clarify the evolutionary origins of environmental ecotypes. By conducting a
284 comparative analysis of genomes within the *Escherichia* cryptic clade IV from diverse
285 ecological (environmental or animal host) and geographical sources, we identified
286 substantial intra-group genomic structuring linked to ecological origins. Our results
287 shed light on the genomic mechanisms underlying the naturalization phenomena within
288 the *Escherichia* genus and include the first genomic characterization of a member of
289 *Escherichia* cryptic clade IV isolated from a highly polluted urban stream in South
290 America. Furthermore, this study provides relevant data for the sanitary management of
291 urban basins.

292     There is an ongoing debate about the ecological niches occupied by the cryptic clades
293     and their phylogenetic relationships within the genus *Escherichia* (Vital *et al.*, 2015;
294     Walk, 2015; Jang *et al.*, 2017). In this context, the cryptic clade IV genome reported here
295     represents a particularly useful new data point, because it is the first isolate of this type
296     obtained from an environmental source in South America, thus expanding the
297     geographical representation of isolates from cryptic clades. Our analysis highlighted
298     consistent relationships between cryptic clades and *E. coli* strains, which is coherent
299     with previous phylogenies based on single nucleotide polymorphisms (SNPs) or other
300     sets of genes (Walk *et al.*, 2009; Walk, 2015; Jang *et al.*, 2017; van der Putten and Mende,
301     2019). In addition, we observed that strains from clades III and IV fall in close
302     phylogenetic proximity. It has been proposed that the latter groups comprise a new
303     species, which was named *E. ruysiae* (van der Putten and Mende, 2019). Clade I strains
304     clustered together with *E. coli* strains, many of which are known to be human
305     pathogens. This observation aligns with previous studies and observational data,
306     suggesting that clade I may represent a versatile taxon with intermediate properties
307     between a free-living lifestyle and enteric commensalism (Chaudhuri and Henderson,
308     2012; Clermont *et al.*, 2013).

309     Regarding genomic features of clinical and epidemiological interest, the genome of
310     339_SF and other clade IV genomes exhibited low abundance of virulence factors and
311     no antibiotic-resistance genes were detected in most of the genomes of the clade. These
312     characteristics, consistent with the findings of other authors on this clade (Ingle *et al.*,
313     2011), have been associated with a free-living lifestyle (Casadevall and Pirofski, 2000;
314     Uddin *et al.*, 2021), as the distribution of these factors is typically linked to specific
315     selection pressures associated with a host-associated life (Jernberg *et al.*, 2010; Becattini *et*
316     *al.*, 2016).

317     Natural environments, such as urban freshwaters, present specific ecological conditions
318     that act as filters for bacterial taxa. Conditions affecting the viability of microorganisms
319     include a wide range of physicochemical characteristics (e.g., temperature, pH and
320     nutrient availability) and stressors (e.g., interspecific competence, solar radiation and
321     xenobiotic compounds) (Vrede, 2005; Jiang and Patel, 2008; Jang *et al.*, 2017). Therefore,
322     the success of bacterial taxa in persisting and proliferating in these environments
323     depends on their physiological aptitude to face these ecological contexts. To identify
324     genomic factors related to this differential fitness, we explored the enrichment of genes
325     within the cryptic clade IV in relation to their origin (environmental or host). This
326     approach allowed the identification of numerous genes likely associated with the free-
327     living ecotype. Many of these genes function in stress response and DNA repair
328     mechanisms. As well, elements associated with the CRISPR defence system were
329     enriched in environmental genomes. These genes likely confer a selective advantage
330     related to the reduction of high viral loads of open natural systems (Jackson *et al.*, 2017).
331     Additionally, genes linked to lipid and carbohydrate metabolic pathways were enriched
332     in the environmental ecotype, suggesting that the utilization of alternative energy
333     sources may be a feature of environmental isolates. Similar deductions were made by
334     Luo et al. (2011), who also found enrichment of factors related to resource acquisition
335     (diol utilization) and survival (lysozyme production; associated with bacterial innate
336     immunity) in environmental strains of *E. coli*. Overall, the abundance of factors with
337     defined actions against environmental stressors or adaptation to distinctive

338 physicochemical characteristics is consistent with free-living ecotypes. Certainly, the
339 limited range of ecological conditions expected in the enteric cavities of homeotherms
340 would not select for such a wide array of genetic features.

341 We did not find differences in genomic features of clinical relevance among clade IV
342 genomes in relation to their origin. In the same vein, neither the phylogenetic analysis
343 nor the gene content-based clustering method completely discriminated clade IV
344 genomes according to their origin. However, our results demonstrate that there are
345 differences in gene content associated with the origin of isolates. These differences
346 could be interpreted as evolutionarily selected traits, reinforcing the idea of ecological
347 niche specialization. Thus, our results highlight a rich genomic diversity within the
348 cryptic clade IV of the *Escherichia* genus. We believe that future work should be
349 oriented to analyze the association of genomic features with phenotypic characteristics
350 (e.g. growth, adhesion) on *Escherichia* isolates from different sources. This kind of
351 studies will further clarify the relationship between genomic features and ecological
352 niches in the *Escherichia* genus.

## Authors contributions

354 **Martín Saraceno:** conceptualization (equal contribution); formal analysis (head);
355 investigation; writing – original draft preparation (head); writing – review & editing
356 (head). **Nicolás Frankel:** formal analysis (supporting); project administration
357 supervision (supporting); writing – review & editing (supporting). **Martín Graziano:**
358 conceptualization (equal contribution); formal analysis (supporting); project
359 administration supervision (head); writing – review & editing (supporting).

## Acknowledgements and conflict of interest statement

## Data availability statement

366 The genome assembly for *Escherichia* cryptic clade IV strain 339_SF has been
367 deposited at DDBJ/ENA/GenBank under the accession JBBUKW000000000
368 (BioProject PRJNA1092102, BioSample SAMN40619928). Accession information and
369 metadata for the genomes of cryptic clade IV and the rest of the strains used are listed in
370 Table S1.

## References

372 Acosta-Dibarrat, J., Enriquez-Gómez, E., Talavera-Rojas, M., Soriano-Vargas, E., Navarro, A.,
373     and Morales-Espinosa, R. (2021) Characterization of commensal Escherichia coli isolates
374     from slaughtered sheep in Mexico. *The Journal of Infection in Developing Countries* **15**:
375     1755–1760.

376 Becattini, S., Taur, Y., and Pamer, E.G. (2016) Antibiotic-Induced Changes in the Intestinal
377     Microbiota and Disease. *Trends Mol Med* **22**: 458–478.

378 Beghain, J., Bridier-Nahmias, A., Nagard, H. Le, Denamur, E., and Clermont, O. (2018)
379     ClermonTyping: An easy-to-use and accurate in silico method for Escherichia genus strain
380     phylotyping. *Microb Genom* **4**: 1–8.

381 Berthe, T., Ratajczak, M., Clermont, O., Denamur, E., and Petit, F. (2013) Evidence for
382     coexistence of distinct escherichia coli populations in various aquatic environments and
383     their survival in estuary water. *Appl Environ Microbiol* **79**: 4684–4693.

384 Bleuel, C., Große, C., Taudte, N., Scherer, J., Wesenberg, D., Krauß, G.J., et al. (2005) TolC is
385     involved in enterobactin efflux across the outer membrane of Escheriia coli. *J Bacteriol*
386     **187**: 6701–6707.

387 Carattoli, A., Zankari, E., Garciá-Fernández, A., Larsen, M.V., Lund, O., Villa, L., et al. (2014)
388     In Silico detection and typing of plasmids using plasmidfinder and plasmid multilocus
389     sequence typing. *Antimicrob Agents Chemother* **58**: 3895–3903.

390 Casadevall, A. and Pirofski, L.A. (2000) Host-pathogen interactions: Basic concepts of
391     microbial commensalism, colonization, infection, and disease. *Infect Immun* **68**: 6511–
392     6518.

393 Chaudhuri, R.R. and Henderson, I.R. (2012) Infection , Genetics and Evolution The evolution of
394     the Escherichia coli phylogeny. *Infection, Genetics and Evolution* **12**: 214–226.

395 Clermont, O., Christenson, J.K., Denamur, E., and Gordon, D.M. (2013) The Clermont
396     Escherichia coli phylo-typing method revisited: Improvement of specificity and detection
397     of new phylo-groups. *Environ Microbiol Rep* **5**: 58–65.

398 Clermont, O., Gordon, D.M., Brisse, S., Walk, S.T., and Denamur, E. (2011) Characterization of
399     the cryptic Escherichia lineages: Rapid identification and prevalence. *Environ Microbiol*
400     **13**: 2468–2477.

401 Duan, Y., Gao, H., Zheng, L., Liu, S., Cao, Y., Zhu, S., et al. (2020) Antibiotic Resistance and
402     Virulence of Extraintestinal Pathogenic Escherichia coli (ExPEC) Vary According to
403     Molecular Types. *Front Microbiol* **11**:.

404 Florensa, A.F., Kaas, R.S., Clausen, P.T.L.C., Aytan-Aktug, D., and Aarestrup, F.M. (2022)
405     ResFinder – an open online resource for identification of antimicrobial resistance genes in
406     next-generation sequencing data and prediction of phenotypes from genotypes. *Microb*
407     *Genom* **8**: 1–10.

408 Gurevich, A., Saveliev, V., Vyahhi, N., and Tesler, G. (2013) QUAST: Quality assessment tool
409     for genome assemblies. *Bioinformatics* **29**: 1072–1075.

410 Harrington, B. (2004) Inkscape.

411 Hildebrand, F., Meyer, A., and Eyre-Walker, A. (2010) Evidence of selection upon genomic GC-
412     content in bacteria. *PLoS Genet* **6**:.

413 Holden, N., Pritchard, L., and Toth, I. (2009) Colonization outwith the colon: Plants as an
414     alternative environmental reservoir for human pathogenic enterobacteria: Review article.
415     *FEMS Microbiol Rev* **33**: 689–703.

416 Ingle, D.J., Clermont, O., Skurnik, D., Denamur, E., Walk, S.T., and Gordon, D.M. (2011)
417     Biofilm formation by and thermal niche and virulence characteristics of Escherichia spp.
418     *Appl Environ Microbiol* **77**: 2695–2700.

419 Ishii, S., Ksoll, W.B., Hicks, R.E., and Sadowsky, M.J. (2006) Presence and growth of
420     naturalized Escherichia coli in temperate soils from lake superior watersheds. *Appl*
421     *Environ Microbiol* **72**: 612–621.

422  Jackson, S.A., McKenzie, R.E., Fagerlund, R.D., Kieper, S.N., Fineran, P.C., and Brouns, S.J.J.
423       (2017) CRISPR-Cas: Adapting to change. *Science (1979)* **356**:.

424  Jang, J., Hur, H.G., Sadowsky, M.J., Byappanahalli, M.N., Yan, T., and Ishii, S. (2017)
425       Environmental Escherichia coli: ecology and public health implications—a review. *J Appl*
426       *Microbiol* **123**: 570–581.

427  Jernberg, C., Löfmark, S., Edlund, C., and Jansson, J.K. (2010) Long-term impacts of antibiotic
428       exposure on the human intestinal microbiota. *Microbiology (N Y)* **156**: 3216–3223.

429  Jiang, L. and Patel, S.N. (2008) Community assembly in the presence of disturbance: A
430       microcosm experiment. *Ecology* **89**: 1931–1940.

431  Lee, C.M., Lin, T.Y., Lin, C.C., Kohbodi, G.N.A., Bhatt, A., Lee, R., and Jay, J.A. (2006)
432       Persistence of fecal indicator bacteria in Santa Monica Bay beach sediments. *Water Res*
433       **40**: 2593–2602.

434  Liu, B., Zheng, D., Zhou, S., Chen, L., and Yang, J. (2022) VFDB 2022: A general classification
435       scheme for bacterial virulence factors. *Nucleic Acids Res* **50**: D912–D917.

436  Luo, C., Walk, S.T., Gordon, D.M., Feldgarden, M., Tiedje, J.M., and Konstantinidis, K.T.
437       (2011) Genome sequencing of environmental Escherichia coli expands understanding of
438       the ecology and speciation of the model bacterial species. *Proceedings of the National*
439       *Academy of Sciences* **108**: 7200–7205.

440  Mackowiak, M., Leifels, M., Hamza, I.A., Jurzik, L., and Wingender, J. (2018) Distribution of
441       Escherichia coli, coliphages and enteric viruses in water, epilithic biofilms and sediments
442       of an urban river in Germany. *Science of The Total Environment* **626**: 650–659.

443  McLellan, S.L., Fisher, J.C., and Newton, R.J. (2015) The microbiome of urban waters.
444       *International Microbiology* **18**: 141–149.

445  Müller, C.M., Åberg, A., Straseviçiene, J., Emody, L., Uhlin, B.E., and Balsalobre, C. (2009)
446       Type 1 fimbriae, a colonization factor of uropathogenic Escherichia coli, are controlled by
447       the metabolic sensor CRP-cAMP. *PLoS Pathog* **5**: e1000303.

448  Nguyen, L.T., Schmidt, H.A., Von Haeseler, A., and Minh, B.Q. (2015) IQ-TREE: A fast and
449       effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol*
450       *Evol* **32**: 268–274.

451  Ozenberger, B.A., Schrodt Nahlik, M., and McIntosh, M.A. (1987) Genetic organization of
452       multiple fep genes encoding ferric enterobactin transport functions in Escherichia coli. *J*
453       *Bacteriol* **169**: 3638–3646.

454  Page, A.J., Cummins, C.A., Hunt, M., Wong, V.K., Reuter, S., Holden, M.T.G., et al. (2015)
455       Roary: Rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**: 3691–3693.

456  van der Putten, B. and Mende, D.R. (2019) Escherichia ruysiae sp. nov., isolated from an
457       international traveller. *Int J Syst Evol Microbiol* **71**: 004609.

458  Py, B., Loiseau, L., and Barras, F. (2001) An inner membrane platform in the type II secretion
459       machinery of Gram-negative bacteria. *EMBO Rep* **2**: 244–248.

460  Rambaut, A. (2018) FigTree.

461  Russo, T.A., Wenderoth, S., Carlino, U.B., Merrick, J.M., and Lesse, A.J. (1998) Identification,
462       genomic organization, and analysis of the group III capsular polysaccharide genes kpsD,

463      kpsM, kpsT, and kpsE from an extraintestinal isolate of Escherichia coli (CP9,
464      O4/K54/H5). *J Bacteriol* **180**: 338–349.

465 Di Sante, L., Pugnaloni, A., Biavasco, F., Giovanetti, E., and Vignaroli, C. (2018) Multicellular
466      behavior of environmental Escherichia coli isolates grown under nutrient-poor and low-
467      temperature conditions. *Microbiol Res* **210**: 43–50.

468 Saraceno, M., Lugo, S.G., Ortiz, N., Gómez, B., y García, C.S., Frankel, N., and Graziano, M.
469      (2020) Unraveling the ecological processes modulating the population structure of
470      Escherichia coli in a highly polluted urban stream network. *Sci Rep* **11**: 1804.

471 Seemann, T. (2020) Abricate.

472 Seemann, T. (2014) Genome analysis Prokka: rapid prokaryotic genome annotation. **30**: 2068–
473      2069.

474 de Sena Brandine, G. and Smith, A.D. (2019) Falco: high-speed FastQC emulation for quality
475      control of sequencing data. *F1000Res* **8**: 1874.

476 Stecher, B. and Hardt, W.D. (2011) Mechanisms controlling pathogen colonization of the gut.
477      *Curr Opin Microbiol* **14**: 82–91.

478 Steinsland, H., Lacher, D.W., Sommerfelt, H., and Whittam, T.S. (2010) Ancestral lineages of
479      human enterotoxigenic Escherichia coli. *J Clin Microbiol* **48**: 2916–2924.

480 Suits, M.D.L., Lang, J., Pal, G.P., Couture, M., and Jia, Z. (2009) Structure and heme binding
481      properties of Escherichia coli O157:H7 ChuX. *Protein Science* **18**: 825–838.

482 Surbeck, C.Q., Jiang, S.C., and Grant, S.B. (2010) Ecological control of fecal indicator bacteria
483      in an urban stream. *Environ Sci Technol* **44**: 631–637.

484 The Consortium Uniprot (2023) UniProt: the Universal Protein Knowledgebase in 2023 -
485      Google Scholar. **51**: 523–531.

486 Toloza, L., Giménez, R., Fábrega, M.J., Alvarez, C.S., Aguilera, L., Cañas, M.A., et al. (2015)
487      The secreted autotransporter toxin (Sat) does not act as a virulence factor in the probiotic
488      Escherichia coli strain Nissle 1917. *BMC Microbiol* **15**: 1–15.

489 Touchon, M., Perrin, A., de Sousa, J.A.M., Vangchhia, B., Burn, S., O'Brien, C.L., et al. (2020)
490      Phylogenetic background and habitat drive the genetic diversification of Escherichia coli.
491      *PLoS Genet* **16**: e1008866.

492 Uddin, T.M., Chakraborty, A.J., Khusro, A., Zidan, B.R.M., Mitra, S., Emran, T. Bin, et al.
493      (2021) Antibiotic resistance in microbes: History, mechanisms, therapeutic strategies and
494      future prospects. *J Infect Public Health* **14**: 1750–1766.

495 Vignaroli, C., Di Sante, L., Magi, G., Luna, G.M., Di Cesare, A., Pasquaroli, S., et al. (2015)
496      Adhesion of marine cryptic Escherichia isolates to human intestinal epithelial cells. *ISME*
497      *Journal* **9**: 508–515.

498 Vital, M., Chai, B., Østman, B., Cole, J., Konstantinidis, K.T., and Tiedje, J.M. (2015) Gene
499      expression analysis of E. coli strains provides insights into the role of gene regulation in
500      diversification. *ISME Journal* **9**: 1130–1140.

501 Vrede, K. (2005) Nutrient and temperature limitation of bacterioplankton growth in temperate
502      lakes. *Microb Ecol* **49**: 245–256.

503 Walk, S.T. (2015) The "Cryptic" Escherichia. *EcoSal Plus* **6**: 10–1128.

504  Walk, S.T., Alm, E.W., Calhoun, L.M., Mladonicky, J.M., and Whittam, T.S. (2007) Genetic
505      diversity and population structure of Escherichia coli isolated from freshwater beaches. **9**:
506      2274–2288.

507  Walk, S.T., Alm, E.W., Gordon, D.M., Ram, J.L., Toranzos, G.A., Tiedje, J.M., and Whittam,
508      T.S. (2009) Cryptic lineages of the genus Escherichia. *Appl Environ Microbiol* **75**: 6534–
509      6544.

510  Wick, R.R., Judd, L.M., Gorrie, C.L., and Holt, K.E. (2017) Unicycler: Resolving bacterial
511      genome assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**: 1–22.

512  Wood, D.E., Lu, J., and Langmead, B. (2019) Improved metagenomic analysis with Kraken 2.
513      *Genome Biol* **20**: 1–13.

514  Zhou, Z., Alikhan, N.F., Mohamed, K., Fan, Y., and Achtman, M. (2020) The EnteroBase user's
515      guide, with case studies on Salmonella transmissions, Yersinia pestis phylogeny, and
516      Escherichia core genomic diversity. *Genome Res* **30**: 138–152.

517