



HAL
open science

The *cnf1* gene is associated to an expanding *Escherichia coli* ST131 H30Rx/C2 sublineage and confers a competitive advantage for host colonization

Landry Laure Tsoumtsa Meda, Luce Landraud, Serena Petracchini, Stéphane Descorps-Declere, Emeline Perthame, Marie-Anne Nahori, Laura Ramirez Finn, Molly A Ingersoll, Rafael Patiño-Navarete, Philippe Glaser, et al.

► To cite this version:

Landry Laure Tsoumtsa Meda, Luce Landraud, Serena Petracchini, Stéphane Descorps-Declere, Emeline Perthame, et al.. The *cnf1* gene is associated to an expanding *Escherichia coli* ST131 H30Rx/C2 sublineage and confers a competitive advantage for host colonization. 2021. pasteur-03797969v1

HAL Id: pasteur-03797969

<https://pasteur.hal.science/pasteur-03797969v1>

Preprint submitted on 15 Oct 2021 (v1), last revised 5 Oct 2022 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

33 **SUMMARY**

34 Epidemiological projections point to acquisition of ever-expanding multidrug resistance
35 (MDR) by *Escherichia coli*, a commensal of the digestive tract acting as a source of urinary
36 tract pathogens. We performed a high-throughput genetic screening of predominantly
37 clinical *E. coli* isolates from wide geographical origins. This revealed a preferential
38 distribution of the Cytotoxic Necrotizing Factor 1 (CNF1)-toxin encoding gene, *cnf1*, in four
39 sequence types encompassing the pandemic *E. coli* MDR lineage ST131. This lineage is
40 responsible for a majority of extraintestinal infections that escape first-line antibiotic
41 treatment and has known enhanced capacities to colonize the gastrointestinal tract (GIT).
42 Statistical modeling uncovered a dominant global expansion of *cnf1*-positive strains within
43 multidrug-resistant ST131 subclade H30Rx/C2. Despite the absence of phylogeographical
44 signals, *cnf1*-positive isolates adopted a clonal distribution into clusters on the ST131-
45 H30Rx/C2 phylogeny, sharing a similar profile of virulence factors and the same *cnf1* allele.
46 Functional analysis of the *cnf1*-positive clinical strain EC131GY ST131-H30Rx/C2, established
47 that a *cnf1*-deleted EC131GY is outcompeted by the wildtype strain in a mouse model of
48 competitive infection of the bladder while both strains behave similarly during
49 monoinfections. This points for positive selection of *cnf1* during UTI rather than
50 urovirulence. Wildtype EC131GY also outcompeted the mutant when concurrently
51 inoculated into the gastrointestinal tract, arguing for selection within the gut. Whatever the
52 site of selection, these findings support that the benefit of *cnf1* enhancing host colonization
53 by ST131-H30Rx/C2 in turn drives a worldwide dissemination of the *cnf1* gene together with
54 extended spectrum of antibiotic resistance genes.

55

56 INTRODUCTION

57 CNF1 is a paradigm of bacterial deamidase toxins activating Rho GTPases¹⁻⁴. Clinical studies
58 document a higher prevalence of the *cnf1*-encoding gene in uropathogenic strains of
59 *Escherichia coli* (UPEC), which belong to the larger group of extraintestinal pathogenic *E. coli*
60 (ExPEC), as compared to commensals from healthy patients⁵⁻⁷. Urinary tract infections (UTI)
61 are common infections that affect more than 150 million individuals annually and are the
62 second cause of antibiotic prescribing⁸. Despite clinical evidence of a role for *cnf1* in
63 urovirulence⁶, attempts to define fitness advantages conferred by this toxin in mouse
64 models of UTI have led to opposing conclusions, although these studies do suggest that
65 CNF1 toxin activity may worsen inflammation and tissue damage⁹⁻¹³. Moreover, in an
66 animal model of bacteremia, CNF1 exerts a paradoxical avirulent effect antagonized by the
67 action of the genetically-associated alpha-hemolysin, further blurring the role of CNF1 in
68 host-pathogen interactions¹⁴⁻¹⁶. In *E. coli*, there are three types of CNF-like toxins sharing
69 high amino acid sequence identities¹⁷⁻²⁰. However, isolates expressing the CNF2 and CNF3
70 toxins are rarely detected in extraintestinal infections in humans. Large-scale population
71 genetics studies to analyse the distribution of *cnf*-like toxin genes in *E. coli* would give
72 important insights regarding their dynamics within the *E. coli* population.

73 *E. coli* represents the predominant aerobic bacteria of the gut microbiota, as well as an
74 extraintestinal opportunistic pathogen^{21,22}. Carriage of ExPEC in the gut is a putative source
75 of extraintestinal infections, including UTIs²³⁻²⁶. Only a few sequence types (STs) within the
76 *E. coli* population account for more than half of all *E. coli* strains responsible for
77 extraintestinal infections not causally related to antibiotic resistance^{21,27}. The globally
78 disseminated *E. coli* ST131 has emerged as the predominant lineage responsible for
79 worldwide dissemination of *bla*_{CTX-M-15} extended spectrum beta-lactamase and the rise of
80 multidrug resistant (MDR) extraintestinal infections^{28,29}. This well-defined clonal group is
81 structured into three different clades, with the fluoroquinolone (FQ)-resistant clade C strains
82 subdivided into two subclades comprised of H30R/C1 and the dominant expanding
83 H30Rx/C2, frequently carrying *bla*_{CTX-M-15}³⁰⁻³². Enhanced interindividual transmission and
84 dispersal of *E. coli* ST131 lineage likely accounts for the lack of phylogeographical signal³³. A
85 larger sampling of strains from the domestic and wild animal world is necessary to better
86 appreciate host specific marks on the evolutionary history of this lineage.

87 One reason for the unprecedented success of *E. coli* ST131-*H30* clade C may be its intrinsic
88 capacity to persist in the gastrointestinal tract (GIT) in competition with other strains of *E.*
89 *coli*^{24,34–37}. Enhanced colonization capacities of the gastrointestinal tract by *E. coli* ST131
90 likely promote inter-individual transmission, favoring its dissemination in the human
91 population and other hosts, as compared to other lineages^{24,38,39}. The remarkable fitness of
92 this lineage strongly supports the idea of a step-wise acquisition of factors promoting gut
93 colonization, potentially scattered in the UPEC populations. In this case, virulence can be
94 considered as a by-product of commensalism, “virulence factors” being in fact selected for
95 increasing fitness in the commensal niche⁴⁰.

96 To better appreciate *cnf1* dynamics, we performed a large-scale screen of the toxin gene
97 distribution in the *E. coli* population. Its increasing prevalence in the ST131-*H30Rx/C2* lineage
98 led us to test whether an advantage is conferred by *cnf1* for GIT colonization. Wildtype
99 EC131GY from ST131-*H30Rx/C2* outcompeted the mutant when concurrently inoculated into
100 the gastrointestinal tract, arguing for selection within the gut. The *cnf1*-deleted EC131GY is
101 also outcompeted by the wildtype strain during competitive infection of the bladder.
102 However, in monoinfections both strains infected similarly, pointing to possible positive
103 selection mechanism for *cnf1* during UTI and demonstrating that *cnf1* is not an urovirulence
104 factor. These findings support that the benefit of *cnf1* enhancing host colonization by
105 ST131-*H30Rx/C2* in turn drives a worldwide dissemination of this lineage.

106

107 RESULTS

108

109 Analysis of the distribution of *cnf* genes in a large collection of *E. coli* genomes

110 At the start of this study, we mined large genomic datasets from EnteroBase to gain more
111 insight into the distribution of the *cnf1* gene and its close homologs in the population of *E.*
112 *coli*⁴¹. EnteroBase represents an integrated software environment widely used to define the
113 population structure of several bacterial genera, including pathogens. Quantitative
114 information on the collection of 141,234 *E. coli* genomes deposited in EnteroBase are
115 reported in the supplementary figure 1. This collection, starting from 1900, aggregates
116 genomes from strains collected worldwide, but mainly in Europe and North America, and
117 from a wide range of sources but principally human isolates (Sup. Figure 1A, 1B, 1C). Using a
118 Hidden Markov Model (HMM) approach, coupled to amino acid pairwise distance
119 calculation, we retrieved *cnf*-like positive strains and characterized each type of *cnf*
120 sequence. In total, we identified $n=6,411$ *cnf*-positive strains (4.5% of all *E. coli* isolates) with
121 a remarkable dominance of *cnf1* (87.8%, $n=5,634$), as compared to *cnf2* (8.6%, $n=554$) and
122 *cnf3* (3.5%, $n=223$). These strains displayed only one CNF-like toxin encoding gene. The
123 prevalent *cnf1* gene in this genomic dataset was widely distributed among isolates of all
124 origins but most notably in the groups denoted humans (5.4% of $n=48,518$ human isolates)
125 and companion animals (24.1% of $n=2,652$ companion animal isolates) (Sup. Figure 1C).

126

127 We next studied the distribution of *cnf1* among *E. coli* phylogenetic groups and sequence
128 types (STs). The *cnf1* gene is preferentially associated with isolates from the phylogroup B2,
129 representing 24.3% of $n=22,305$ retrieved genome sequences (Sup. Figure 1D). We observed
130 a tight association of *cnf1* with the most frequently encountered ExPEC sequence types (STs)
131 (Table 1). Notably, a majority of the 5,634 *cnf1*-positive strains segregated among the four
132 sequence types: ST131 (24.5% of *cnf1*-positive strains, $n=1,382$), ST73 (23.2%, $n=1,308$), ST12
133 (12.4%, $n=699$) and ST127 (10.7%, $n=601$) with the remaining 29.2% of *cnf1*-positive strains
134 widely distributed among 266 other STs. Interestingly, we noticed a steady increase of the
135 percentage of *cnf1*-positive strains in the *E. coli* ST131 lineage from 13% in 2009 up to 23% in
136 2019 (Figure 1), while this percentage fluctuated around high values in ST73, ST12 and
137 ST127. This analysis reveals a close association of *cnf1* with common ExPEC lineages and a
138 surprising convergent distribution of *cnf1* in ST131, ST73 and ST127 that are representative

139 of adherent-invasive *E. coli* (AIEC) associated with colonic Crohn's disease and known to
140 have enhanced capacities to colonize the gastrointestinal tract^{21,42,43}.

141

142 ***cnf1*-positive strains segregate into monophyletic groups in ST131 phylogeny**

143 The rising prevalence of *cnf1* in *E. coli* ST131 motivated us to study its distribution in this
144 lineage, as its phylogenetic structure is well defined and displays a major FQ-resistant clade
145 largely independent of geographical signal^{30–33}. EnteroBase contained 9,242 genomes of *E.*
146 *coli* ST131 at the time of analysis (November 2020). To ease genomic analysis, we retained
147 5,231 genomes that were isolated from 1967 to 2018. We built a Maximum Likelihood
148 phylogenetic tree based on a total of 37,304 non-recombinant SNPs. Phylogenetic
149 distribution of strains showed an expected dominant population of clade C (76%, $n = 3,981$;
150 99% *fimH30*), as compared to clade A (11%, $n = 569$; 92% *fimH41*) and B (13%, $n = 68$; 62%
151 *fimH22*) (Figure 2A, detailed in Sup. Figure 2A). We also found an expected co-distribution of
152 *parC* (S80I/E84V) and *gyrA* (S83L/D87N) alleles that confer strong resistance to FQ in most
153 strains from clade C (99.84%, $n=3,975$ strains), and a tight association of the *bla*_{CTX-M-15} ESBL
154 gene (85%, $n=2,194$ isolates) with strains from subclade *H30Rx/C2*. The high number of
155 strains gave enough resolution to distinguish two sublineages, C2_1 and C2_2, originating
156 from C2_0 (Figure 2A). From available metadata, we verified the absence of overall
157 geographical and temporal links in the phylogenetic distribution of *E. coli* ST131 strains (Sup.
158 Figure 2B). In conclusion, large scale phylogenetic reconstruction of ST131 genomes from
159 EnteroBase showed an expected phylogenetic distribution within clades and subclades of
160 genetic traits defining this lineage.

161

162 We next analyzed the distribution of *cnf1*-positive strains ($n=725$) in *E. coli* ST131 phylogeny
163 (Figure 2A, black stripes). The *cnf1*-positive strains were preferentially associated with clade
164 C2 ($n=520$), as compared to clade C1 ($n=101$), clade B ($n=72$) and clade A ($n=32$) (Figure 2A).
165 Strikingly, most *cnf1*-positive strains segregated into lineages in all clades and subclades with
166 a noticeable distribution of *cnf1*-positive ST131 strains in two large lineages (LL) in *H30R/C1*
167 ($n=101$ *cnf1*-positive strains/107 strains in CNF1_LL1) and in *H30Rx/C2_1* ($n=396$ *cnf1*-
168 positive strains/425 strains in the CNF1_LL2) (Figure 2A). We then analyzed the diversity of
169 alleles of *cnf1* to define their distribution in ST131 phylogeny (Sup. Table 1). A similar
170 analysis was performed with the alpha-hemolysin encoding gene, *hlyA*. We found a wide co-

171 distribution of one combination of alleles of *cnf1* (allele P1_{*cnf1*}, 85,1%) and alpha-hemolysin
172 encoding gene *hlyA* (allele P1_{*hlyA*}, 77,2%) in *E. coli* ST131 clade A and C, whereas strains from
173 clade B displayed a large range of combinations of various alleles (Sup. Figure 2A). Together,
174 our data point to a clonal expansion of worldwide disseminated ST131-*H30* strains having
175 the same allele of *cnf1*. Together, this prompted us to perform a clustering analysis of ST131-
176 *H30* strains according to their accessory gene contents. We generated a pan-genome matrix
177 of 51,742 coding sequences from the $n=3,981$ strains of clade C. The dataset of accessory
178 genes was built from $n=7,678$ sequences that were present in at least 50 and no more than
179 3,931 strains. We conducted a hierarchical clustering of strains according to the Ward's
180 minimum variance-derived method⁴⁴ and retained 10 distinct accessory gene clusters.
181 Strikingly, this revealed a conservation between phylogenetically-defined groups CNF1_LL1
182 and CNF1_LL2 and groups defined by their accessory gene contents (Figure 2B). Indeed, the
183 hierarchical clustering was most evident for CNF1_LL2, showing a differential enrichment of
184 $n=1,434$ genes as compared to other strains from clade C, determined with Scoary
185 (Bonferroni-adjusted P -value <0.05)⁴⁵. Together, these data point towards intensive group-
186 specific diversification of accessory gene content in *cnf1*-positive clusters in ST131-*H30*.

187

188 ***cnf1*-positive strains of *E. coli* ST131 segregate between two clade-specific virulence** 189 **profiles**

190 We then defined strain contents in virulence factors (VF) and acquired antibiotic-resistance
191 genes (RG) to perform an unbiased analysis of their distribution into clusters, using a latent
192 block model approach. Acquired antibiotic-resistance genes in ST131 genomes were
193 identified with ResFinder⁴⁶. Profiles of virulence factors were defined with the database
194 published by Petty and colleagues³¹. The unsupervised clustering procedure retained a total
195 of 10 RG-clusters and 7 VF-clusters (Figure 3A). Differences in number of VFs and RGs
196 between clusters were all significant (Figure 3B). We found that *cnf1*-positive strains were
197 scattered among several RG clusters (Figure 3A, left panel). By contrast, most *cnf1*-positive
198 strains segregated into the cluster VF4 (84% of *cnf1*-positive strains, $n=609$) with the
199 remaining 16% strains being distributed between VF1 (15%) and other VF clusters (1%)
200 (Figure 3A, right panel). In contrast to RG-clusters, we observed that VF-clusters formed
201 phylogenetically defined groups (Figure 3C). A majority of *cnf1*-positive strains from clade A
202 and B were positive for the VF1 cluster, whereas *cnf1*-positive strains from clade C were

203 positive for the VF4 cluster. With a mean value of 33 virulence factors (Figure 3B), VF4-
204 positive strains displayed the largest arsenal of virulence factors. The VF1 profile was more
205 specifically defined by the presence of genes encoding the IbeA invasin and IroN Salmochelin
206 siderophore receptor (Sup. Figure 3A). By contrast, major determinants of the VF4 cluster
207 encompassed *cnf1* and *hlyA* (54% and 61% in VF4 versus 34% in VF1 and 3% in all other VFs).
208 Specific VF determinants of VF4 also encompassed genes encoding the UclD adhesin that
209 tipped F17-like chaperone-usher (CU) fimbriae cluster and PapG II adhesin from
210 pyelonephritis-associated pili (*pap*) operon (Sup. Figure 3A)^{47,48}. These elements can be
211 genetically associated and constitute the backbone of *cnf1*-bearing pathogenicity islands
212 (PAI) II_{J96} from the O4:K6 *E. coli* strain J96, although PAI II_{J96} contains a *papG* class III
213 sequence (Sup. Figure 3B). In good agreement, analysis of several complete sequences of
214 *cnf1*-bearing PAI II_{J96}-like from ST131-*H30* showed a conservation of a module containing this
215 set of genes, defining VF4 (Sup. Figure 3B).

216

217 ***cnf1*-positive strains display dominant expansion in ST131-*H30Rx/C2***

218 We next analyzed the temporal distribution of *cnf1*-positive strains within clades and
219 subclades. Using a Generalized Linear Models (GLM) approach, we first verified within our
220 dataset the increase of *fimH30*-positive isolates over time (clade C) in *E. coli* ST131 that was
221 maximal in *H30Rx/C2* ($P < 2 \cdot 10^{-16}$) (Figure 4A). We also noted a significant increase in the
222 proportion of *cnf1*-positive strains over time in *E. coli* ST131 (Figure 4B, top panel). The GLM
223 was then fitted on years, clades, and subclades. We tested the significance of the year effect
224 and *P*-values were corrected for multiple comparisons using Tukey's method. The year effect
225 was not significant for clade A, B, or subclade *H30R/C1* (Figure 4B). Instead, we observed a
226 significant increase of the proportion of *cnf1*-positive strains within *H30Rx/C2* over time
227 ($P = 1.25 \cdot 10^{-11}$). In addition, the GLM fitted curves predicted that the prevalence of *cnf1*-
228 positive strains within *H30Rx/C2* sublineage would be approximately 50% (confidence
229 interval of 95% [43% to 58%] in 2018; [47% to 64%] in 2019). Predictive values were
230 confronted to the prevalence of *cnf1* in ST131 strains isolated in 2018 or 2019 in a second
231 independent dataset up-loaded from Enterobase in September 2020. This confirmed the
232 rising prevalence of *cnf1*-positive strains within the sublineage *H30Rx/C2* up to 45% in 2018
233 and 48% in 2019. In conclusion, we identified a dominant expansion of *cnf1*-positive strains
234 within ST131-*H30Rx/C2*.

235

236 ***cnf1* confers a competitive advantage for bladder infection and gut colonization in a ST131-**
237 **H30Rx/C2 strain**

238 The dominant expansion of *cnf1*-positive strains in ST131 H30Rx/C2 prompted us to explore
239 whether CNF1 confers a competitive advantage for bladder infection and/or intestinal
240 colonization. In the cohort SEPTICOLI of bloodstream infections in human adults ⁴⁹, we
241 identified a VF4/*cnf1*-positive strain of *E. coli* ST131 H30Rx/C2, here referred to as EC131GY
242 (Sup. Figure 4). This strain is amenable to genetic engineering and displays a *cnf1*-bearing PAI
243 (PAI II_{EC131GY}) highly similar to the prototypic PAI II_{J96} from the J96 (O4:H5:K6) UPEC strain
244 (Sup. Figure 3B) ⁵⁰. We generated a EC131GY strain in which *cnf1* was replaced with a
245 kanamycin resistance cassette (EC131GYΔ*cnf1*::*kan^r*) and verified the absence of CNF1
246 expression (Sup. Figure 5A). We next verified, *in vitro*, the absence of fitness cost due to the
247 kanamycin resistance cassette as shown by equal growth of parental and Δ*cnf1*::*kan^r*
248 EC131GY strains, and the absence of competition between the strains when grown together
249 (Sup. Figure 5B and 5C). Considering the tight association of *cnf1* with clinical strains of *E. coli*
250 responsible for UTI, we first investigated the impact of the toxin during concurrent infection
251 of the bladder with wild-type EC131GY and EC131GYΔ*cnf1*::*kan^r*. Wild-type *E. coli*
252 outcompeted the isogenic *cnf1*-deficient EC131GY in the first 24 hours, when bacteria must
253 rapidly establish their niche in the face of passive and innate immune host defenses (Figure
254 5A). This fitness advantage was maintained at day 3 and 7, demonstrating that *cnf1* plays a
255 role in the early stages of UPEC pathogenesis, as previously suggested ⁹. No difference of
256 colonization of wild-type EC131GY and EC131GYΔ*cnf1*::*kan^r* was observed in monomicrobial
257 bladder infections (Figure 5B). This finding can be interpreted as a positive selection
258 mechanism to maintain the CNF1 gene during UTI, considering that a loss of *cnf1* would be
259 detrimental for bacterial fitness in a mixed population. We then explored the impact of *cnf1*
260 in GIT colonization, again by competitive infection with EC131GY WT and
261 EC131GYΔ*cnf1*::*kan^r*, using intra-gastric gavage ⁵¹. Longitudinal measurements of CFU in the
262 feces showed that CNF1 conferred an advantage to wild-type EC131GY over the
263 EC131GYΔ*cnf1*::*kan^r* isogenic strain for gut colonization from 9 days after oral gavage, which
264 persisted over 27 days (Figure 5B). Together, these data uncover the advantage conferred by
265 CNF1 in a setting of competitive UTI and for intestinal colonization by the VF4/*cnf1*-positive
266 EC131GY strain from the ST131-H30Rx/C2 lineage.

267 **DISCUSSION**

268 Initially thought to be absent in the *Escherichia coli* ST131 lineage, the *cnf1* gene was
269 estimated to be found in approximately 15% of this lineage, among 99 isolates from distinct
270 geographical locations across the world, in 2014^{31,52}. Large-scale genetic analysis of more
271 than five thousand isolates of *E. coli* ST131 from Enterobase, a database widely used by
272 clinicians, gives here sufficient statistical power to unveil a dominant expansion trend of
273 *cnf1*-positive strains within clade H30Rx/C2. Our analysis supports the hypothesis of a recent
274 expansion of a large phylogenetic subcluster of *cnf1*-positive ST131-H30Rx/C2 strains
275 circulating between humans and dogs^{53,54}. In addition, we document a stable population
276 dynamic of *cnf1*-positive H30R/C1 strains within clade C1. This raises the question of
277 whether *cnf1* confers a fitness advantage at the population level. Our compelling findings
278 ascribed such a feature of *cnf1* to specific genetic backgrounds, thereby enhancing the
279 expansion and dissemination of a subpopulation of ST131-H30Rx/C2 within the ST131
280 lineage. Furthermore, we report the high prevalence of *cnf1* gene in the three sequence
281 types ST73, ST12 and ST127 of *E. coli* that have different antibiotic resistance profiles.
282 Together, this points to a role of *cnf1* in the dynamics of ExPEC that is independent from
283 antibiotic resistance genetic backgrounds. The rising prevalence of *cnf1*-positive H30Rx/C2,
284 and evidence of their mobilization between humans and dogs⁵³, suggest that *cnf1* enhances
285 the dissemination of H30Rx/C2 within households with companion animals, which is likely
286 driven by an increased ability to compete for GIT colonization. In further support of this
287 conclusion, we found a prevalence of 24% of *cnf1*-positive strains in the group companion
288 animals from the Enterobase database. Finally, we report a high occurrence of the *cnf1* gene
289 in common AIEC pathotypes responsible for Crohn's disease and known to colonize the GIT
290 well^{21,42,43}. These findings highlight the importance of studying the interplay between CNF1
291 and the gut mucosa for persistence and inflammatory bowel diseases.

292 The competitive advantage conferred by *cnf1* during the acute phase of UTI (i.e., 24 hours)
293 suggests this toxin promotes FimH-dependent invasion of urothelial cells, which results in
294 the formation of intracellular bacterial communities (IBCs)^{8,55}. In support of this hypothesis,
295 cell biology studies show that CNF1 promotes invasion of host cells by *E. coli* through its
296 capacity to activate host Rho GTPases^{20,56-58}. Although this remains to be formally
297 demonstrated, CNF1 deamidase likely exacerbates the activation of Rho GTPases, which are
298 required for type I pili-mediated host cell invasion⁵⁹. Importantly, in contrast to concurrent

299 infection, *cnf1* confers no detectable virulence advantage during bladder mono-infection.
300 Considering that UTI caused by *E. coli* are usually dominated by one strain, we propose that
301 the fitness advantage conferred by *cnf1* during concurrent infection could reflect a positive
302 selection mechanism to maintain the gene during UTI. Alternatively, as *cnf1* also confers a
303 fitness advantage in the gut commensal niche which is the primary *E. coli* habitat, the
304 selective pressure occurs in the gut and *cnf1* confers virulence as a by-product of
305 commensalism⁴⁰. This mechanism has been shown for the PAIs of the B2 ST127 strain 536⁶⁰.
306 The F17-like pilus adhesin UclD from *cnf1*-bearing PAI confers a competition advantage for
307 gut colonization, while it shows no virulence role in UTI⁵¹. Therefore, this also points for
308 *cnf1*-driven positive selection as a potential broader mechanism to maintain the PAI during
309 UTI.

310 Our findings that *cnf1* gives a competitive advantage for GIT colonization also raise the
311 interest of defining epistatic relationships between factors encoded within the core set of
312 genes of the PAI II_{EC131GY} from ST131 H30Rx/C2 for colonization and bacterial persistence in
313 tissues. Indeed, these operons encode F17-like pili, the P-fimbriae tipped with PapG class II
314 adhesin, and the *hlyA* toxin, as well as a gene encoding haemagglutinin in *E. coli* K1 (Hek)^{61–}
315⁶³. This also includes elements of oxidative stress adaptation, namely the methionine
316 sulfoxide reductase complex MsrPQ encoding genes *yedYZ*, which may work against CNF1-
317 generated oxidative stress^{64,65}.

318 Collectively, our findings point towards a bidirectional interplay between *cnf1* and the *E. coli*
319 ST131 lineage to enhance host colonization by H30Rx/C2 whatever the site of selection and
320 to promote a worldwide dissemination of the Cytotoxic Necrotizing Factor 1-encoding gene
321 together with extended spectrum of antibiotic resistant genes.

322

323 **FIGURE LEGENDS**

324 **Figure 1: Prevalence overtime in representative *E. coli* sequence types bearing *cnf1***

325 Bar chart show number of *E. coli* strains from ST131, ST127, ST73 and ST12 isolated each
326 year during the period 2002-2019, left y-axis. Percentages of *cnf1*-positive strains per year,
327 right y-axis.

328

329 **Figure 2: Dynamic of CNF1-encoding gene in *E. coli* ST131 from EnteroBase**

330 **A)** Maximum likelihood phylogeny of *E. coli* ST131 from EnteroBase (Sup. Figure 2 for
331 extended information). The phylogeny was constructed with 5,231 genomes for a total of
332 37,304 non-recombinant core-genome SNPs. The different clades and subclades A, B, C0, C1,
333 C2_0, C2_1, C2_2 are highlighted in blue, red, light green, green, pink, orange and purple
334 respectively. From inside to outside circles are indicated (1) *fimH* alleles, (2) *gyrA* and *parC*
335 alleles conferring resistance to FQ (shown in green), (3) strains positive for *bla*_{CTX-M-15} (shown
336 in orange) and (4) strains bearing *cnf1* gene (shown in black). **B)** Hierarchical clustering of
337 strains from clade C (*n* = 3981 strains) based on their accessory gene content. The pan-
338 genome is composed of 51,742 genes including 2,672 genes that are present in 98% of the
339 strains. The graph displays the 7,678 genes identified as present in at least 50 and less than
340 3,930 genomes. The colored annotation indicates (from left to right) the presence of *cnf1*
341 (CNF1_status), clades (C1, C1 CNF1_LL1, C2_0, C2_1, C2_1 CNF1_LL2, C2_2) and accessory
342 genes cluster (AG_clusters). Large lineages of *cnf1*-positive strains in clades C1 and C2_1 are
343 denoted CNF1_LL1 and CNF1_LL2, respectively.

344

345 **Figure 3: Co-clustering of acquired antibiotic-resistance gene and virulence factors in *E. coli***
346 **ST131.**

347 **A)** Heatmaps show clusters of antibiotic acquired-resistance gene (RG) (left panel) or
348 virulence gene (VF) (right panel) profiles (Sup. table 2) constructed using a binary latent
349 block model between strains by row and RGs or VFs by column. Black lines indicate the
350 presence of RG or VF in each strain. Annotations are displayed on the right of each heatmap:
351 information about strain clusters and *fimH* alleles together with *hlyA* and *cnf1* carriage. **B)**
352 Box-and-whisker plot showing the distribution of strains according to their content of
353 acquired antibiotic-resistance genes (upper panel) or content of virulence factors (lower
354 panel). The dotted line shows the mean number of RG or VF. All one-versus-all comparisons

355 of VF and RG contents between clusters ($*P < 0.05$, $***P < 0.001$). **C)** RG, VF clusters and
356 *cnf1* carriage are displayed on the *E. coli* ST131 phylogenetic tree. The different clades and
357 subclades A, B, C0, C1, C2_0, C2_1, C2_2 are highlighted in blue, red, light green, green, pink,
358 orange and purple respectively.

359

360 **Figure 4: Increase over the year in the proportion of *cnf1*-positive strains in *E. coli* ST131**
361 **H30Rx/C2**

362 **A)** Distribution of *fimH* alleles (upper panel) or clades/subclades (lower panel) within the
363 study population of *E. coli* ST131. Both figures show observed counts per year (dots) and
364 data fitted lines (dashed lines) with a generalized linear model (Poisson regression). **B)**
365 Increase of the proportion of *cnf1*-positive strains in the whole *E. coli* ST131 population
366 along time (top panel, $P = 7.41 \cdot 10^{-7}$) and by clades and subclades. The black dots represent
367 the observed proportion of *cnf1*-positive strains by year with fitted line of a logistic
368 regression model (blue curves). Dashed grey lines display the 95% confidence intervals. The
369 *P*-values are not significant for clade A ($P = 0.287$), B ($P = 0.952$), H30R/C1 ($P = 0.992$) and
370 significant for H30Rx/C2 ($P = 1.25 \cdot 10^{-11}$).

371

372 **Figure 5: CNF1 promotes ST131-H30Rx/C2 bladder and intestinal colonization**

373 Mice were infected concurrently **(A)** or separately **(B)** with wild-type EC131GY (WT) and
374 EC131GY $\Delta cnf1::kan^r$ ($\Delta cnf1$) via intravesical instillation of the bladder. For GIT colonization,
375 mice were pretreated with streptomycin and subsequently infected concurrently via the oral
376 route with EC131GY WT and $\Delta cnf1$ **(C)**. Levels of viable bacteria in bladder homogenates or
377 feces were assessed at indicated times by measuring colony forming units (CFU). Data
378 represent the competitive index (CI) (A and C) or CFU per bladder (B) for each animal and
379 medians (red bar). Total of $n=15-18$ (bladder CI, three replicates), $n=9-10$ (bladder single,
380 two replicates at day 1) and $n=21$ (intestine, three replicates). $*P < 0.05$, $**P < 0.01$, $***P <$
381 0.001 , $****P < 0.0001$ and ns : non-significant by Wilcoxon signed-rank test.

382

383 **Table 1: Distribution of phylogroups and sequence types among *E. coli* *cnf*-positive strains**
384 **from EnteroBase**

385 The total number and the percentage of each phylogroup and most dominant sequence
386 types (STs) among *cnf*-positive strains are indicated

387 **MATERIAL and METHODS**

388 ***E. coli* genome collection**

389 Collection of 141,234 *E. coli* genome sequences from EnteroBase (November 2020)
390 (<http://enterobase.warwick.ac.uk>)⁴¹. Strain's metadata (collection year, continent, source
391 niche of isolation and sequence type) were also retrieved (Sup. Table 3). Assemblies were
392 downloaded in GenBank format and proteomes generated using annotations provided in
393 GenBank files.

394

395 ***In silico* detection and typing of CNF-like toxin encoding genes**

396 The search for *cnf* genes in *E. coli* genomes was carried out with a domain specific Hidden
397 Markov Models (HMM) profile built with 16 representative sequences of CNF1 catalytic
398 domain (Sup. Table 4) using HMMER (<http://hmmer.org/>)⁶⁶. Protein sequences from
399 positive hits were extracted from EnteroBase annotated *E. coli* proteomes and submitted to
400 Clustal Omega for the computation of pairwise distances of the sequences, along with
401 representative sequences of CNF-like toxin (CNF1 (AAA85196.1), CNF2 (WP_012775889.1)
402 and CNF3 (WP_02231387.1)). Distances were used to determine the type of toxin with a
403 threshold value of 0.1. In total 2.7% of HMM-positive sequences with a threshold value
404 above 0.1 against all type of CNF-like toxin or below 0.1 against at least two type of CNF-like
405 toxin were excluded from the analysis.

406

407 **ST131 dataset structure and phylogenomic analysis**

408 The database used for phylogenetic and statistical analyses consists of whole-genome
409 sequences of *E. coli* ST131 isolates collected by mining EnteroBase from 1967 to 2018⁴¹.
410 Leaning on Find ST(s) tool from EnteroBase, we retained a total of 5,231 genome assemblies
411 and associated metadata, including information of the isolation date, country and source of
412 isolates (Sup. table 5). Phylogeny of ST131 isolates was resolved using core non-recombinant
413 SNPs defined with Parsnp (in total 37,304 SNPs)⁶⁷ and Gubbins v2.3.4⁶⁸. A maximum-
414 likelihood tree was then estimated with RAxML v8.2.8 applying a general time-reversible
415 substitution-model with a gamma distribution rate across sites and with an ascertainment
416 bias correction⁶⁹ and the resulting tree was edited with the interactive Tree of Life (iTol) v4
417 program⁷⁰.

418

419 ***In silico* antimicrobial resistance and virulence-associated markers**

420 GyrA and ParC protein sequences were retrieved from the EnteroBase annotated genomes,
421 and aligned with the mafft L-INS-I approach ⁷¹. After a visual inspection of the alignment, in-
422 house customized perl scripts (<https://github.com/rpatinonavarrete/QRDR>) were used to
423 identify the amino acids at the quinolone resistance-determining region (QRDR) (positions
424 83 and 87, and 80 and 84 in GyrA and ParC, respectively). Search for *cnf1* and *hlyA* alleles in
425 ST131 genomes dataset was carried out by Blastn analysis. Sequences were next aligned
426 with Muscle ⁷² and curated to remove incomplete sequences. SNPs were then extracted
427 using SNP-sites ⁷³. To determine strain specific VF profiles, annotated VFs from UPEC
428 described in ³¹ were translated and pBLASTed against ST131 genomes dataset considering
429 only hits with e-value < 10⁻⁵ and identical matches > 95% (sup. Table 2) ⁷⁴. Acquired
430 antibiotic-resistance genes (RGs) in ST131 genomes were defined with ResFinder ⁴⁶.

431

432 **Generalized linear model**

433 Proportion of *cnf1* along time was modeled using a generalized linear model (logistic
434 regression) adjusted on the effect of years and clades with an interaction between these two
435 factors. First, to test if the evolution of *cnf1* proportion was either specific to each clade or
436 global, the significance of the interaction term was tested with a likelihood ratio test, which
437 compares the above-mentioned model against the null model, with no interaction. Then, we
438 investigated the possible increase of the proportion of *cnf1* within each clade. The
439 significance of the slope coefficient for each clade was tested by computing contrasts of the
440 above model. *P*-values were adjusted for multiplicity using single-step correction method.
441 The distribution of *fimH* alleles and clades/subclades within the study population of *E. coli*
442 ST131 was analyzed with a similar approach, except that a Poisson regression model was
443 used to model counting data. The hypothesis testing strategy to investigate the significance
444 of the increase of *fimH* alleles and clades/subclades along time is discussed above.

445

446 **Co-clustering method**

447 Statistical analyses were performed using R software version 3.6.0. A total of 20 strains from
448 the collection of 5,231 strains of *E. coli* ST131 were removed from the analysis due to
449 incomplete associated metadata. The clustering of strains with specific virulence or acquired
450 antibiotic-resistance gene profiles was performed with binary latent block model,

451 implemented in the R package blockcluster ⁷⁵. In this package, the model, a mixture of
452 Bernoulli distributions proposed by ⁷⁶, is estimated using an efficient EM algorithm. As
453 proposed by the authors, the number of clusters was estimated by maximizing the ICL
454 criterion on a bidimensional grid of parameters making this unsupervised classification
455 procedure automatic.

456

457 **Pan-genome analysis**

458 The pangenome of *E. coli* ST131 was estimated using Roary, a high-speed pan genome
459 pipeline analysis tool ⁷⁷. Roary returns as output, the gene presence/absence matrix. The
460 matrix was curated to retain genes present in at least 50 genomes and less than 3980
461 genomes (7678 sequences), that constituted our accessory genes pool dataset. Hierarchical
462 clustering analysis was then conducted by using the pheatmap package in R ([cran.r-
463 project.org/web/packages/pheatmap/index.html](http://cran.r-project.org/web/packages/pheatmap/index.html)). The gene presence/absence file
464 generated by Roary was further analyzed using Scoary ⁴⁵ with a significant Bonferroni-
465 adjusted P-value < 0.05 for genes associated to *cnf1*-positive lineages (Sup. Table 8).

466

467 **Mouse colonization model**

468 Local Animal Studies Committee and National Research Council approved all procedures
469 used for the mouse experiments described in the present study (APAFIS#26133-
470 202006221228936 v1, 2016–0010. For gut colonization, groups of female C57BL/6 mice aged
471 6–7 weeks (Charles River) were pretreated with a single dose of streptomycin (1 g/kg in 200
472 μ l water) *per os* 1 day prior to gavage, as described in ⁵¹. The strains derived from the clinical
473 strain H1-001-0141-G-Y, here referred to as EC131GY (de Lastours et al., 2020), are
474 described in the extended materials and methods section. Mice were co-infected *per os* with
475 2×10^9 CFU of each strain in 200 μ l PBS. Fecal pellets were collected from every individual
476 mouse at indicated times, weighed and homogenized in 500 μ l phosphate-buffered saline
477 (PBS) pH 7.2 by vigorous vortexing. CFUs were determined by plating serial dilutions on
478 selective LB agar plates. Strains were prepared for infection as follows: a single colony of
479 EC131GY or its derivative was inoculated in 10 ml selective LB medium and incubated at 37°C
480 under static conditions for 24h. Bacteria were then inoculated in 25 ml fresh selective LB
481 medium at 1:1000 dilution and incubated at 37°C under static conditions for 18-24h.
482 Bacteria were then washed twice in cold PBS, and concentrated in PBS at approximately

483 2×10^9 CFU per 200 μ l. Inocula titers are verified in parallel for each infection. For intravesical
484 infection: Urinary tract infection was induced in mice as previously described^{78,79}. Briefly, a
485 single colony of EC131GY or the *cnf1* mutant was inoculated in 10 ml LB medium with
486 antibiotics and incubated at 37°C under static conditions for 18h. Mice were infected with a
487 total of 10^7 CFU of bacteria in 50 μ l PBS via a rigid urinary catheter under anesthesia. To
488 calculate CFU, bladders were aseptically removed and homogenized in 1 ml of PBS. Serial
489 dilutions were plated on LB agar plates with antibiotics, as required. The competitive index
490 (CI) was calculated as: CFU WT output strain/CFU mutant output strain, with the verification
491 in each experiment that CFU WT input strain/CFU mutant input strain was close to 1. A
492 Wilcoxon signed-rank test was performed to assess the statistical significance of differences
493 in CI over time. Statistical analyses were performed using GraphPad Prism 9.

494

495

496 **ACKNOWLEDGMENTS**

497 This work was supported by the “Fondation ARC” PJA 20191209650, the “Fondation pour la
498 Recherche Médicale” (Equipe FRM 2016, DEQ20161136698), Ligue Nationale contre le
499 Cancer Subvention de Recherche Scientifique, RS20/75-63 and the French National Research
500 Agency (ANR-10-LABX-62-IBEID, INCEPTION) and ANR-17-CE17-0014. The plasmid pKOBEG
501 was kindly provided by Jean-Marc Ghigo.

502

503 **AUTHOR CONTRIBUTIONS**

504 Bioinformatics analyses were performed L.T.M., S.D.-D., R.P.N. and analyzed by E.L., L.L., P.G.
505 and E.D. Statistical analyses were performed by L.L. and E.P. *In vivo* experiments were
506 coordinated by A.M., M.A.I., O.D. and performed by M.-A. N., A.M. and L.R.F. with strains
507 engineered by S.P. and A.M. The research was coordinated by E.L. and manuscript drafted
508 with help of L.T.M., L.L., O.D., E.D. and P.G. Manuscript was reviewed and approved by all
509 authors.

510

511 **REFERENCES**

512

513 1. Flatau, G. et al. Toxin-induced activation of the G protein p21 Rho by deamidation of
514 glutamine. *Nature* **387**, 729-733 (1997).

515 2. Schmidt, G. et al. Gln 63 of Rho is deamidated by *Escherichia coli* cytotoxic necrotizing factor-1.
516 *Nature* **387**, 725-729 (1997).

517 3. Aktories, K. & Barbieri, J. T. Bacterial cytotoxins: targeting eukaryotic switches. *Nat Rev*
518 *Microbiol* **3**, 397-410 (2005).

519 4. Patel, J. C. & Galan, J. E. Manipulation of the host actin cytoskeleton by Salmonella--all in the
520 name of entry. *Curr Opin Microbiol* **8**, 10-15 (2005).

521 5. Landraud, L., Gauthier, M., Fosse, T. & Boquet, P. Frequency of *Escherichia coli* strains
522 producing the cytotoxic necrotizing factor (CNF1) in nosocomial urinary tract infections. *Lett*
523 *Appl Microbiol* **30**, 213-216 (2000).

524 6. Dubois, D. et al. Cyclomodulins in urosepsis strains of *Escherichia coli*. *J Clin Microbiol* **48**, 2122-
525 2129 (2010).

526 7. Starčič Erjavec, M. & Žgur-Bertok, D. Virulence potential for extraintestinal infections among
527 commensal *Escherichia coli* isolated from healthy humans--the Trojan horse within our gut.
528 *FEMS Microbiol Lett* **362**, (2015).

529 8. Klein, R. D. & Hultgren, S. J. Urinary tract infections: microbial pathogenesis, host-pathogen
530 interactions and new treatment strategies. *Nat Rev Microbiol* **18**, 211-226 (2020).

531 9. Rippere-Lampe, K. E., O'Brien, A. D., Conran, R. & Lockman, H. A. Mutation of the gene
532 encoding cytotoxic necrotizing factor type 1 (*cnf1*) attenuates the virulence of uropathogenic
533 *Escherichia coli*. *Infect Immun* **69**, 3954-3964 (2001).

534 10. Rippere-Lampe, K. E. et al. Cytotoxic necrotizing factor type 1-positive *Escherichia coli* causes
535 increased inflammation and tissue damage to the prostate in a rat prostatitis model. *Infect*
536 *Immun* **69**, 6515-6519 (2001).

537 11. Garcia, T. A., Ventura, C. L., Smith, M. A., Merrell, D. S. & O'Brien, A. D. Cytotoxic necrotizing
538 factor 1 and hemolysin from uropathogenic *Escherichia coli* elicit different host responses in
539 the murine bladder. *Infect Immun* **81**, 99-109 (2013).

540 12. Michaud, J. E., Kim, K. S., Harty, W., Kasprenski, M. & Wang, M. H. Cytotoxic Necrotizing Factor-
541 1 (CNF1) does not promote *E. coli* infection in a murine model of ascending pyelonephritis.
542 *BMC Microbiol* **17**, 127 (2017).

543 13. Schreiber, H. L. et al. Bacterial virulence phenotypes of *Escherichia coli* and host susceptibility
544 determine risk for urinary tract infections. *Sci Transl Med* **9**, (2017).

- 545 14. Landraud, L., Gibert, M., Popoff, M. R., Boquet, P. & Gauthier, M. Expression of *cnf1* by
546 *Escherichia coli* J96 involves a large upstream DNA region including the hlyCABD operon, and is
547 regulated by the RfaH protein. *Mol Microbiol* **47**, 1653-1667 (2003).
- 548 15. Diabate, M. et al. *Escherichia coli* alpha-Hemolysin Counteracts the Anti-Virulence Innate
549 Immune Response Triggered by the Rho GTPase Activating Toxin CNF1 during Bacteremia. *PLoS*
550 *Pathog* **11**, e1004732 (2015).
- 551 16. Dufies, O. et al. *Escherichia coli* Rho GTPase-activating toxin CNF1 mediates NLRP3
552 inflammasome activation via p21-activated kinases-1/2 during bacteraemia in mice. *Nat*
553 *Microbiol* **6**, 401-412 (2021).
- 554 17. Falbo, V., Pace, T., Picci, L., Pizzi, E. & Caprioli, A. Isolation and nucleotide sequence of the gene
555 encoding cytotoxic necrotizing factor 1 of *Escherichia coli*. *Infect Immun* **61**, 4909-4914 (1993).
- 556 18. Orden, J. A. et al. Necrotoxicogenic *Escherichia coli* from sheep and goats produce a new type of
557 cytotoxic necrotizing factor (CNF3) associated with the *eae* and *ehxA* genes. *Int Microbiol* **10**,
558 47-55 (2007).
- 559 19. Oswald, E. et al. Cytotoxic necrotizing factor type 2 produced by virulent *Escherichia coli*
560 modifies the small GTP-binding proteins Rho involved in assembly of actin stress fibers. *Proc*
561 *Natl Acad Sci U S A* **91**, 3814-3818 (1994).
- 562 20. Ho, M., Mettouchi, A., Wilson, B. A. & Lemichez, E. CNF1-like deamidase domains: common
563 Lego bricks among cancer-promoting immunomodulatory bacterial virulence factors. *Pathog*
564 *Dis* **76**, doi: 10.1093/femspd/fty045 (2018).
- 565 21. Denamur, E., Clermont, O., Bonacorsi, S. & Gordon, D. The population genetics of pathogenic
566 *Escherichia coli*. *Nat Rev Microbiol* **19**, 37-54 (2021).
- 567 22. Tenaillon, O., Skurnik, D., Picard, B. & Denamur, E. The population genetics of commensal
568 *Escherichia coli*. *Nat Rev Microbiol* **8**, 207-217 (2010).
- 569 23. Nielsen, K. L., Dynesen, P., Larsen, P. & Frimodt-Møller, N. Faecal *Escherichia coli* from patients
570 with *E. coli* urinary tract infection and healthy controls who have never had a urinary tract
571 infection. *J Med Microbiol* **63**, 582-589 (2014).
- 572 24. Johnson, J. R. et al. Household Clustering of *Escherichia coli* Sequence Type 131 Clinical and
573 Fecal Isolates According to Whole Genome Sequence Analysis. *Open Forum Infect Dis* **3**,
574 ofw129 (2016).
- 575 25. Yamamoto, S. et al. Genetic evidence supporting the fecal-perineal-urethral hypothesis in
576 cystitis caused by *Escherichia coli*. *J Urol* **157**, 1127-1129 (1997).
- 577 26. Moreno, E. et al. Relationship between *Escherichia coli* strains causing acute cystitis in women
578 and the fecal *E. coli* population of the host. *J Clin Microbiol* **46**, 2529-2534 (2008).

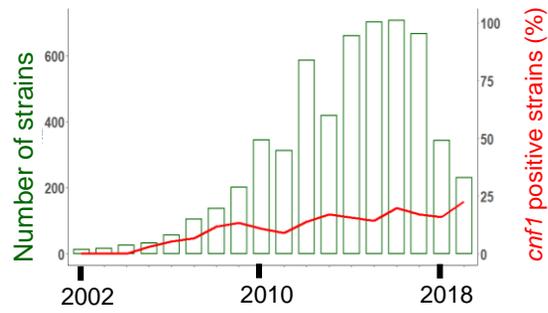
- 579 27. Kallonen, T. et al. Systematic longitudinal survey of invasive *Escherichia coli* in England
580 demonstrates a stable population structure only transiently disturbed by the emergence of
581 ST131. *Genome Res* **27**, 1437-1449 (2017).
- 582 28. Johnson, J. R., Johnston, B., Clabots, C., Kuskowski, M. A. & Castanheira, M. *Escherichia coli*
583 sequence type ST131 as the major cause of serious multidrug-resistant *E. coli* infections in the
584 United States. *Clin Infect Dis* **51**, 286-294 (2010).
- 585 29. Peirano, G. & Pitout, J. D. Molecular epidemiology of *Escherichia coli* producing CTX-M beta-
586 lactamases: the worldwide emergence of clone ST131 O25:H4. *Int J Antimicrob Agents* **35**, 316-
587 321 (2010).
- 588 30. Price, L. B. et al. The epidemic of extended-spectrum- β -lactamase-producing *Escherichia coli*
589 ST131 is driven by a single highly pathogenic subclone, H30-Rx. *MBio* **4**, e00377-13 (2013).
- 590 31. Petty, N. K. et al. Global dissemination of a multidrug resistant *Escherichia coli* clone. *Proc Natl*
591 *Acad Sci U S A* **111**, 5694-5699 (2014).
- 592 32. Ben Zakour, N. L. et al. Sequential Acquisition of Virulence and Fluoroquinolone Resistance Has
593 Shaped the Evolution of *Escherichia coli* ST131. *MBio* **7**, e00347-16 (2016).
- 594 33. McNally, A. et al. Combined Analysis of Variation in Core, Accessory and Regulatory Genome
595 Regions Provides a Super-Resolution View into the Evolution of Bacterial Populations. *PLoS*
596 *Genet* **12**, e1006280 (2016).
- 597 34. Madigan, T. et al. Extensive Household Outbreak of Urinary Tract Infection and Intestinal
598 Colonization due to Extended-Spectrum β -Lactamase-Producing *Escherichia coli* Sequence Type
599 131. *Clin Infect Dis* **61**, e5-12 (2015).
- 600 35. Tchesnokova, V. L. et al. Pandemic Uropathogenic Fluoroquinolone-resistant *Escherichia coli*
601 Have Enhanced Ability to Persist in the Gut and Cause Bacteriuria in Healthy Women. *Clin Infect*
602 *Dis* **70**, 937-939 (2020).
- 603 36. Shevchenko, S. G., Radey, M., Tchesnokova, V., Kisiela, D. & Sokurenko, E. V. *Escherichia coli*
604 Clonobiome: Assessing the Strain Diversity in Feces and Urine by Deep Amplicon Sequencing.
605 *Appl Environ Microbiol* **85**, (2019).
- 606 37. Vimont, S. et al. The CTX-M-15-producing *Escherichia coli* clone O25b: H4-ST131 has high
607 intestine colonization and urinary tract infection abilities. *PLoS One* **7**, e46547 (2012).
- 608 38. Gurnee, E. A. et al. Gut Colonization of Healthy Children and Their Mothers With Pathogenic
609 Ciprofloxacin-Resistant *Escherichia coli*. *J Infect Dis* **212**, 1862-1868 (2015).
- 610 39. Laupland, K. B., Church, D. L., Vidakovich, J., Mucenski, M. & Pitout, J. D. Community-onset
611 extended-spectrum beta-lactamase (ESBL) producing *Escherichia coli*: importance of
612 international travel. *J Infect* **57**, 441-448 (2008).

- 613 40. Le Gall, T. et al. Extraintestinal virulence is a coincidental by-product of commensalism in B2
614 phylogenetic group *Escherichia coli* strains. *Mol Biol Evol* **24**, 2373-2384 (2007).
- 615 41. Zhou, Z. et al. The Enterobase user's guide, with case studies on Salmonella transmissions,
616 *Yersinia pestis* phylogeny, and *Escherichia coli* core genomic diversity. *Genome Res* **30**, 138-152
617 (2020).
- 618 42. Mirsepasi-Lauridsen, H. C. et al. Secretion of Alpha-Hemolysin by *Escherichia coli* Disrupts Tight
619 Junctions in Ulcerative Colitis Patients. *Clin Transl Gastroenterol* **7**, e149 (2016).
- 620 43. Boudeau, J., Glasser, A. L., Masseret, E., Joly, B. & Darfeuille-Michaud, A. Invasive ability of an
621 *Escherichia coli* strain isolated from the ileal mucosa of a patient with Crohn's disease. *Infect*
622 *Immun* **67**, 4499-4509 (1999).
- 623 44. Murtagh, F. & Legendre, P. Ward's Hierarchical Agglomerative Clustering Method: Which
624 Algorithms Implement Ward's Criterion? *Journal of Classification* **31**, 274-295 (2004).
- 625 45. Brynildsrud, O., Bohlin, J., Scheffer, L. & Eldholm, V. Rapid scoring of genes in microbial pan-
626 genome-wide association studies with Scoary. *Genome Biol* **17**, 238 (2016).
- 627 46. Zankari, E. et al. Identification of acquired antimicrobial resistance genes. *J Antimicrob*
628 *Chemother* **67**, 2640-2644 (2012).
- 629 47. Blum, G., Falbo, V., Caprioli, A. & Hacker, J. Gene clusters encoding the cytotoxic necrotizing
630 factor type 1, Prs-fimbriae and alpha-hemolysin form the pathogenicity island II of the
631 uropathogenic *Escherichia coli* strain J96. *FEMS Microbiol Lett* **126**, 189-195 (1995).
- 632 48. Bidet, P. et al. Multiple insertional events, restricted by the genetic background, have led to
633 acquisition of pathogenicity island I/J96-like domains among *Escherichia coli* strains of different
634 clinical origins. *Infect Immun* **73**, 4081-4087 (2005).
- 635 49. de Lastours, V. et al. Mortality in *Escherichia coli* bloodstream infections: antibiotic resistance
636 still does not make it. *J Antimicrob Chemother* **75**, 2334-2343 (2020).
- 637 50. Swenson, D. L., Bukanov, N. O., Berg, D. E. & Welch, R. A. Two pathogenicity islands in
638 uropathogenic *Escherichia coli* J96: cosmid cloning and sample sequencing. *Infect Immun* **64**,
639 3736-3743 (1996).
- 640 51. Spaulding, C. N. et al. Selective depletion of uropathogenic *E. coli* from the gut by a FimH
641 antagonist. *Nature* **546**, 528-532 (2017).
- 642 52. Nicolas-Chanoine, M. H. et al. Intercontinental emergence of *Escherichia coli* clone O25:H4-
643 ST131 producing CTX-M-15. *J Antimicrob Chemother* **61**, 273-281 (2008).
- 644 53. Bonnet, R. et al. Host Colonization as a Major Evolutionary Force Favoring the Diversity and the
645 Emergence of the Worldwide Multidrug-Resistant *Escherichia coli* ST131. *mBio* **12**, e0145121
646 (2021).

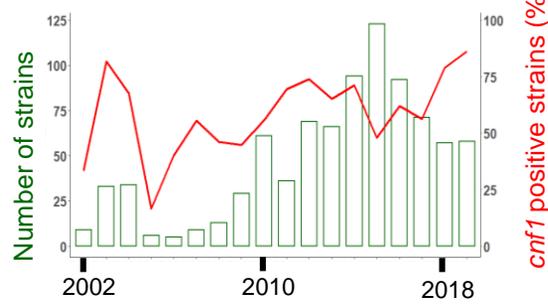
- 647 54. Royer, G. et al. Phylogroup stability contrasts with high within sequence type complex
648 dynamics of *Escherichia coli* bloodstream infection isolates over a 12-year period. *Genome Med*
649 **13**, 77 (2021).
- 650 55. Mulvey, M. A. et al. Induction and evasion of host defenses by type 1-piliated uropathogenic
651 *Escherichia coli*. *Science* **282**, 1494-1497 (1998).
- 652 56. Falzano, L., Rivabene, R., Fabbri, A. & Fiorentini, C. Epithelial cells challenged with a Rac-
653 activating *E. coli* cytotoxin acquire features of professional phagocytes. *Toxicol In Vitro* **16**, 421-
654 425 (2002).
- 655 57. Doye, A. et al. CNF1 exploits the ubiquitin-proteasome machinery to restrict Rho GTPase
656 activation for bacterial host cell invasion. *Cell* **111**, 553-564 (2002).
- 657 58. Visvikis, O. et al. *Escherichia coli* Producing CNF1 Toxin Hijacks Tollip to Trigger Rac1-
658 Dependent Cell Invasion. *Traffic* **12**, 579-590 (2011).
- 659 59. Martinez, J. J. & Hultgren, S. J. Requirement of Rho-family GTPases in the invasion of Type 1-
660 piliated uropathogenic *Escherichia coli*. *Cell Microbiol* **4**, 19-28 (2002).
- 661 60. Tourret, J., Diard, M., Garry, L., Matic, I. & Denamur, E. Effects of single and multiple
662 pathogenicity island deletions on uropathogenic *Escherichia coli* strain 536 intrinsic extra-
663 intestinal virulence. *Int J Med Microbiol* **300**, 435-439 (2010).
- 664 61. Fagan, R. P. & Smith, S. G. The Hek outer membrane protein of *Escherichia coli* is an auto-
665 aggregating adhesin and invasins. *FEMS Microbiol Lett* **269**, 248-255 (2007).
- 666 62. Ristow, L. C. & Welch, R. A. RTX Toxins Ambush Immunity's First Cellular Responders. *Toxins*
667 (*Basel*) **11**, (2019).
- 668 63. Geibel, S. & Waksman, G. The molecular dissection of the chaperone-usher pathway. *Biochim*
669 *Biophys Acta* **1843**, 1559-1567 (2014).
- 670 64. Gennaris, A. et al. Repairing oxidized proteins in the bacterial envelope using respiratory chain
671 electrons. *Nature* **528**, 409-412 (2015).
- 672 65. Falzano, L., Rivabene, R., Santini, M. T., Fabbri, A. & Fiorentini, C. An *Escherichia coli* cytotoxin
673 increases superoxide anion generation via rac in epithelial cells. *Biochem Biophys Res Commun*
674 **283**, 1026-1030 (2001).
- 675 66. Mistry, J., Finn, R. D., Eddy, S. R., Bateman, A. & Punta, M. Challenges in homology search:
676 HMMER3 and convergent evolution of coiled-coil regions. *Nucleic Acids Res* **41**, e121 (2013).
- 677 67. Treangen, T. J., Ondov, B. D., Koren, S. & Phillippy, A. M. The Harvest suite for rapid core-
678 genome alignment and visualization of thousands of intraspecific microbial genomes. *Genome*
679 *Biol* **15**, 524 (2014).
- 680 68. Croucher, N. J. et al. Rapid phylogenetic analysis of large samples of recombinant bacterial
681 whole genome sequences using Gubbins. *Nucleic Acids Res* **43**, e15 (2015).

- 682 69. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
683 phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).
- 684 70. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new developments.
685 *Nucleic Acids Res* **47**, W256-W259 (2019).
- 686 71. Katoh, K., Kuma, K., Miyata, T. & Toh, H. Improvement in the accuracy of multiple sequence
687 alignment program MAFFT. *Genome Inform* **16**, 22-33 (2005).
- 688 72. Edgar, R. C. MUSCLE: a multiple sequence alignment method with reduced time and space
689 complexity. *BMC Bioinformatics* **5**, 113 (2004).
- 690 73. Page, A. J. et al. SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments.
691 *Microb Genom* **2**, e000056 (2016).
- 692 74. Camacho, C. et al. BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
- 693 75. Bhatia, P. S., Iovleff S & Govaert G. blockcluster: An R Package for Model-Based Co-Clustering.
694 *Journal of Statistical Software* **76**, 1-24 (2017).
- 695 76. Govaert, G. & Nadif, M. Block clustering with Bernoulli mixture models: Comparison of
696 different approaches. *Computational Statistics & Data Analysis* **52**, 3233-3245 (2008).
- 697 77. Sitto, F. & Battistuzzi, F. U. Estimating Pangenomes with Roary. *Mol Biol Evol* **37**, 933-939
698 (2020).
- 699 78. Mora-Bau, G. et al. Macrophages Subvert Adaptive Immunity to Urinary Tract Infection. *PLoS*
700 *Pathog* **11**, e1005044 (2015).
- 701 79. Zychlinsky Scharff, A., Albert, M. L. & Ingersoll, M. A. Urinary Tract Infection in a Small Animal
702 Model: Transurethral Catheterization of Male and Female Mice. *J Vis Exp* **130**, 54432 (2017).
- 703

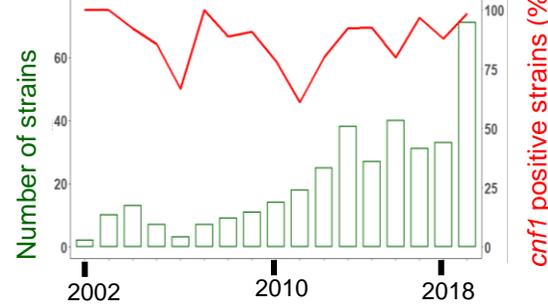
ST131



ST73



ST12



ST127

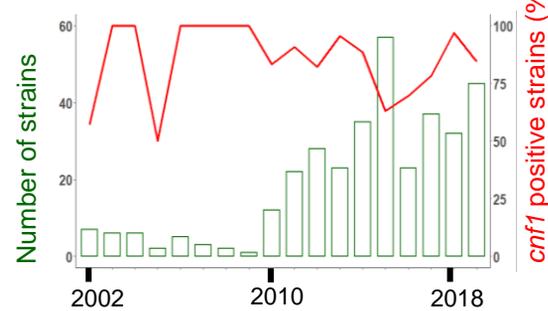


Fig. 1

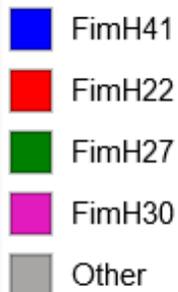
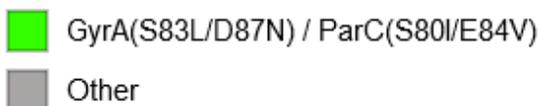
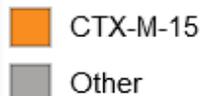
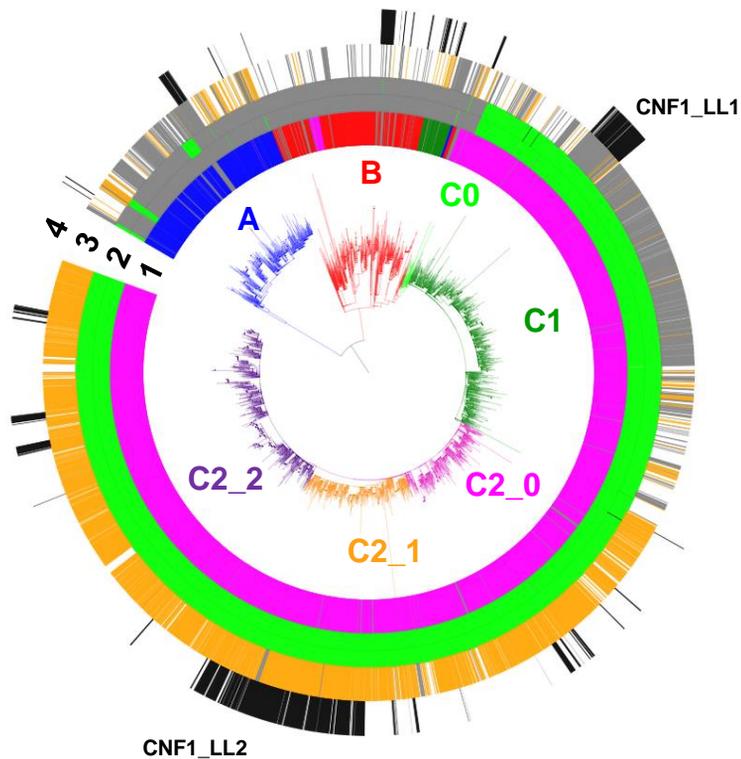
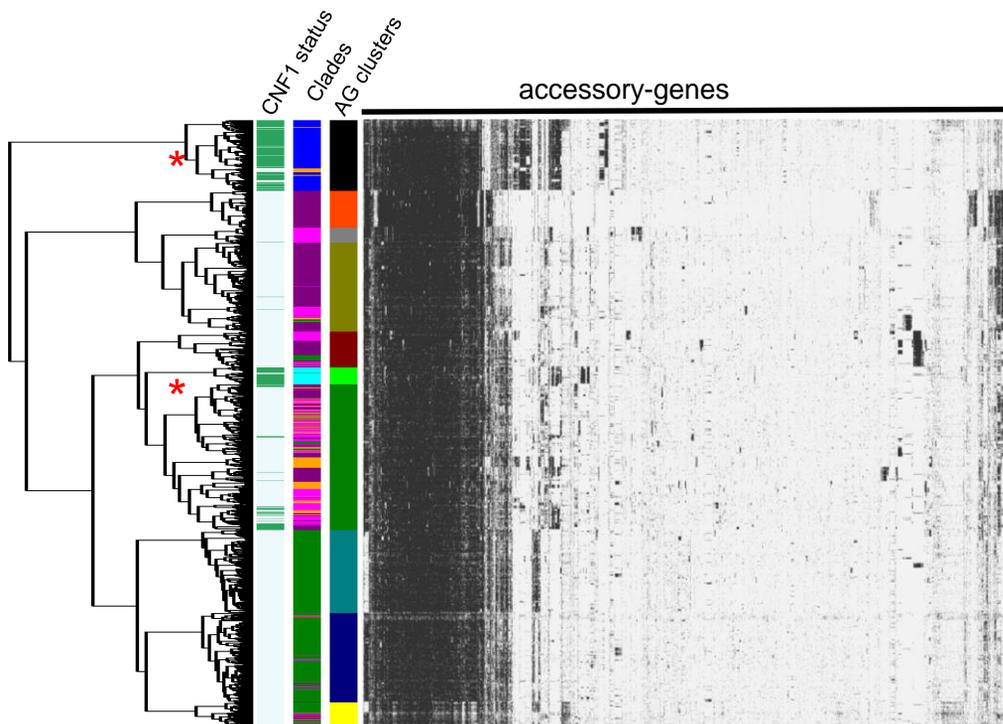
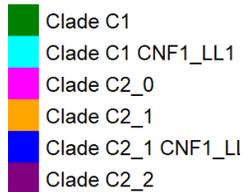
A**1-FimH****2-GyrA/ParC QRDR****3-CTX-M****4-CNF1****B****CNF1_status****Clades**

Fig. 2

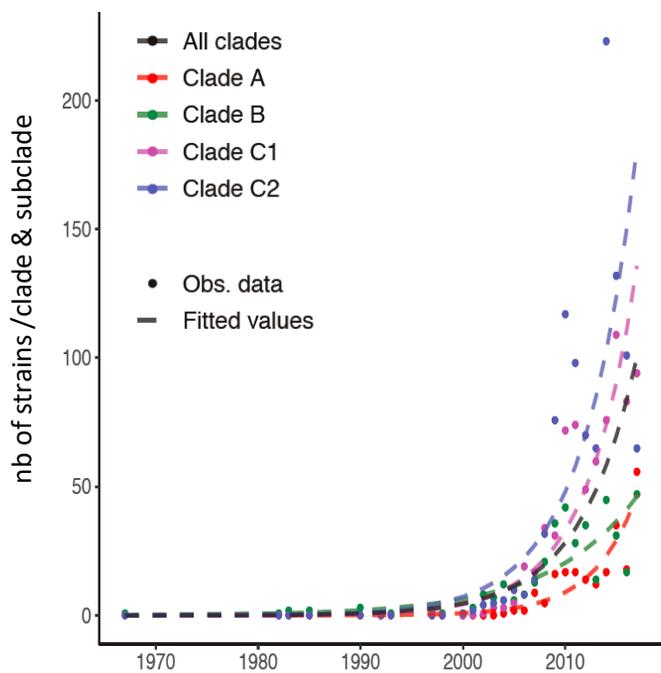
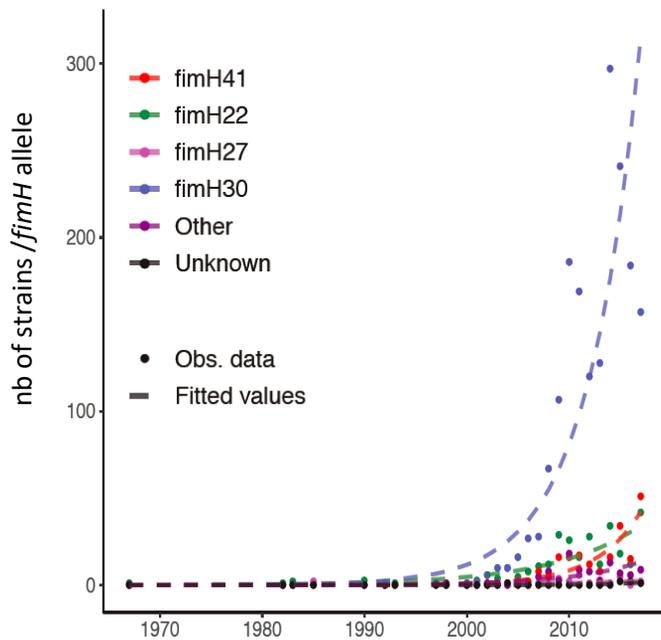
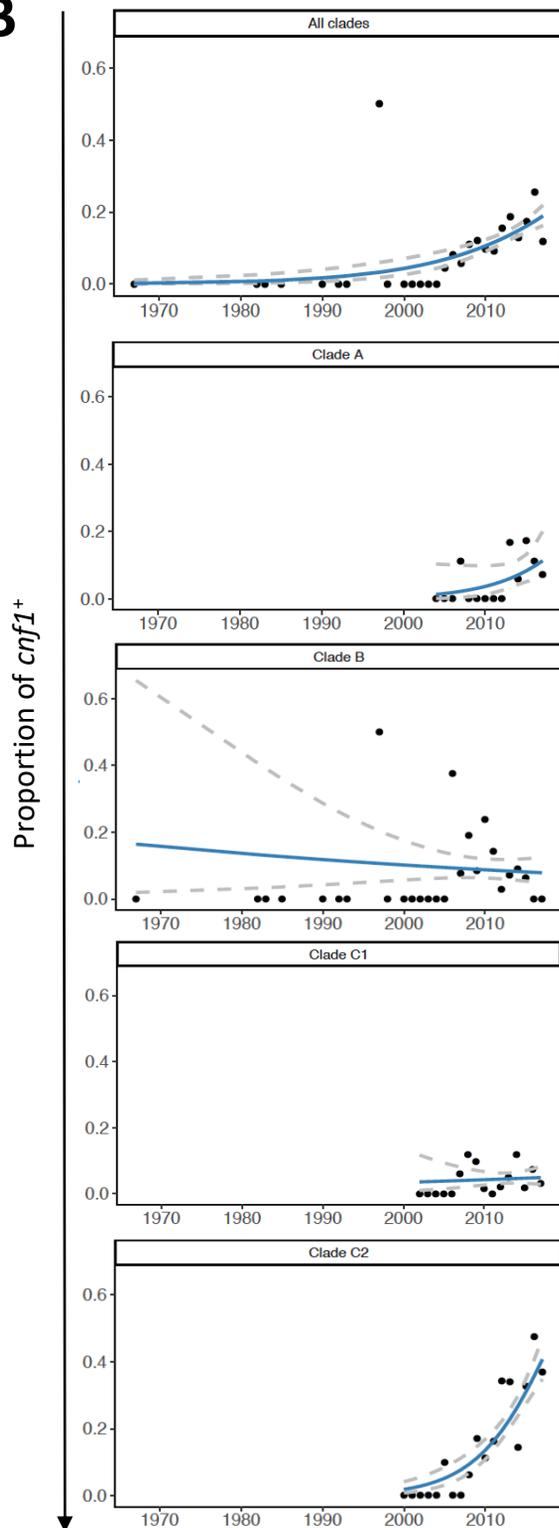
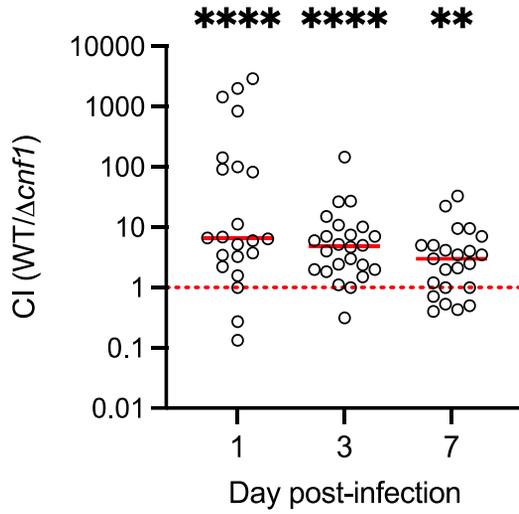
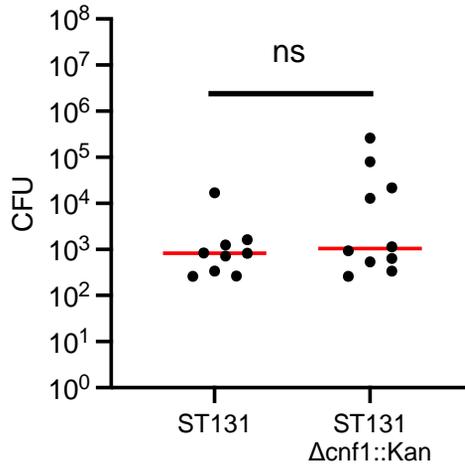
A**B**

Fig. 4

A Competitive index - Bladder



B Single infection - Bladder



C Competitive index - Feces

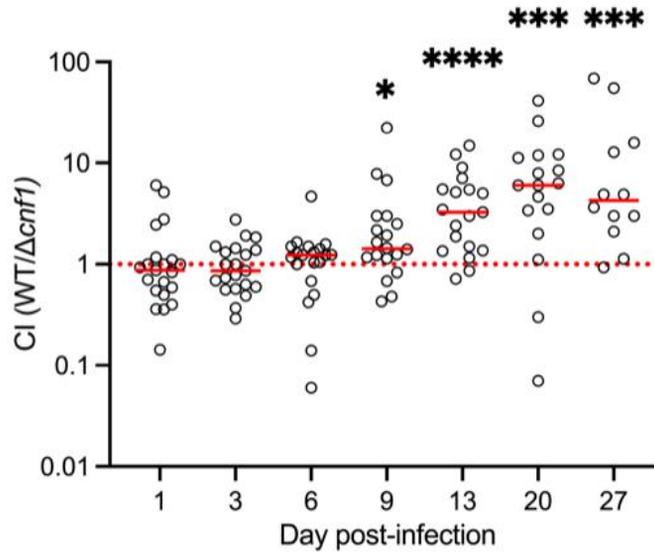


Fig. 5

Phylogroups	ST	Number of strains					Percentage of Phylogroup or Sequence type in CNF-positive strains		
		All	CNF+	CNF1+	CNF2+	CNF3+	CNF1	CNF2	CNF3
A	Total A	34,982	51	0	28	23	0	5.05	10.31
	ST10	8,748	24	0	17	7	0.0	3.1	3.1
	ST342	325	16	0	0	16	0.0	0.0	7.2
B1	Total B1	37,262	527	96	373	58	1.7	67.3	26.0
	ST101	938	93	24	69	0	0.4	12.5	0.0
	ST392	79	66	0	66	0	0.0	11.9	0.0
	ST58	1,487	44	9	35	0	0.2	6.3	0.0
	ST29	496	35	0	0	35	0.0	0.0	15.7
	ST2217	46	31	0	31	0	0.0	5.6	0.0
	ST5738	24	23	0	23	0	0.0	4.2	0.0
	ST21	5,082	10	0	0	10	0.0	0.0	4.5
	ST343	134	2	0	0	2	0.0	0.0	0.9
	ST2836	63	2	0	0	2	0.0	0.0	0.9
ST4063	3	2	0	0	2	0.0	0.0	0.9	
B2	Total B2	22,305	5,478	5,414	63	1	96.1	11.4	0.4
	ST131	9,242	1,383	1,382	0	1	24.5	0.0	0.4
	ST73	2,071	1,308	1,308	0	0	23.2	0.0	0.0
	ST12	809	699	699	0	0	12.4	0.0	0.0
	ST127	709	601	601	0	0	10.7	0.0	0.0
	ST372	366	206	206	0	0	3.7	0.0	0.0
	ST95	1,882	173	147	26	0	2.6	4.7	0.0
	ST141	360	164	164	0	0	2.9	0.0	0.0
	ST998	175	149	149	0	0	2.6	0.0	0.0
	ST80	152	109	105	4	0	1.9	0.7	0.0
	ST537	50	35	35	0	0	0.6	0.0	0.0
ST647	28	26	0	26	0	0.0	4.7	0.0	
C	Total C	3,465	56	45	10	1	0.8	1.8	0.4
D	Total D	9,905	37	20	13	4	0.4	2.3	1.8
E	Total E	16,391	155	7	14	134	0.1	2.5	60.1
	ST11	13,639	113	0	0	113	0.0	0.0	50.7
	ST5592	5	5	0	0	5	0.0	0.0	2.2
ST11457	4	4	0	0	4	0.0	0.0	1.8	
F	Total F	2,957	38	37	0	1	0.7	0.0	0.4
G	Total G	1,862	34	0	34	0	0.0	6.1	0.0
	ST117	1,383	31	0	31	0	0.0	5.6	0.0
Clade I	Total CI	406	18	0	18	0	0.0	3.2	0.0
ST3057	41	11	0	11	0	0.0	2.0	0.0	
Clade II	Total CII	6	0	0	0	0	0.0	0.0	0.0
Clade III	Total CIII	39	0	0	0	0	0.0	0.0	0.0
Clade IV	Total CIV	39	0	0	0	0	0.0	0.0	0.0
Clade V	Total CV	166	0	0	0	0	0.0	0.0	0.0
	Other 358 STs	34,599	1,044	803	215	26	14.3	38.8	11.7

Table 1.