

1 **Insights into the acquisition of the *pks* island and production of colibactin in the**  
2 ***Escherichia coli* population**

3

4

5 Frédéric Auvray<sup>§1</sup>, Alexandre Perrat<sup>1</sup>, Yoko Arimizu<sup>3</sup>, Camille V. Chagneau<sup>1</sup>, Nadège

6 Bossuet-Greif<sup>1</sup>, Clémence Massip<sup>1,2</sup>, Hubert Brugère<sup>1</sup>, Jean-Philippe Nougayrède<sup>1</sup>, Tetsuya

7 Hayashi<sup>3</sup>, Priscilla Branchu<sup>1</sup>, Yoshitoshi Ogura<sup>4</sup>, Eric Oswald<sup>§1,2</sup>

8

9 <sup>1</sup>IRSD, INSERM, Université de Toulouse, INRA, ENVT, UPS, Toulouse, France

10 <sup>2</sup>CHU Toulouse, Hôpital Purpan, Service de Bactériologie-Hygiène, Toulouse, France

11 <sup>3</sup> Department of Bacteriology, Kyushu University, Fukuoka, Japan

12 <sup>4</sup> Division of Microbiology, Department of Infectious Medicine, Kurume University School  
13 of Medicine, Kurume, Fukuoka, Japan

14

15 § Corresponding authors: frederic.auvray@envt.fr; eric.oswald@inserm.fr

16

17 **Keywords:**

18 *pks*, pathogenicity island, genetic diversity, colibactin, genotoxin, *Escherichia coli*,  
19 enterobacteria.

20

21

22 **ABSTRACT**

23 The *pks* island codes for the enzymes necessary for synthesis of the genotoxin colibactin, which  
24 contributes to the virulence of *Escherichia coli* strains and is suspected of promoting colorectal  
25 cancer. From a collection of 785 human and bovine *E. coli* isolates, we identified 109 strains  
26 carrying a highly conserved *pks* island, mostly from the phylogroup B2, but also from  
27 phylogroups A, B1 and D. Different scenarios of *pks* acquisition were deduced from whole  
28 genome sequence and phylogenetic analysis. In the main scenario, *pks* was introduced and  
29 stabilized into certain sequence types (ST) of the B2 phylogroup, such as ST73 and ST95, at  
30 the *asnW* tRNA locus located in the vicinity of the yersiniabactin-encoding High Pathogenicity  
31 Island (HPI). In a few B2 strains, *pks* inserted at the *asnU* or *asnV* tRNA loci close to the HPI  
32 and occasionally was located next to the remnant of an integrative and conjugative element. In  
33 a last scenario specific to B1/A strains, *pks* was acquired, independently of the HPI, at a non-  
34 tRNA locus. All the *pks*-positive strains except 18 produced colibactin. Sixteen strains  
35 contained mutations in *clbB* or *clbD*, or a fusion of *clbJ* and *clbK* and were no longer genotoxic  
36 but most of them still produced low amount of potentially active metabolites associated with  
37 the *pks* island. One strain was fully metabolically inactive without *pks* alteration, but colibactin  
38 production was restored by overexpressing the ClbR regulator. In conclusion, the *pks* island is  
39 not restricted to human pathogenic B2 strains and is more widely distributed in the *E. coli*  
40 population, while preserving its functionality.

41

42

## 43 **IMPACT STATEMENT**

44 Colibactin, a genotoxin associated with the carcinogenicity of certain strains of *E. coli*, is  
45 encoded by a pathogenicity island called *pks*. We took advantage of a large collection of non-  
46 clinical *E. coli* strains originating from human and bovine hosts to explore the distribution,  
47 conservation and functionality of the *pks* island. We found that the *pks* island was not only  
48 present in the phylogroup B2 (and more specifically to certain B2 sublineages), but also in  
49 other genetic phylogroups, highlighting its capacity to disseminate through horizontal gene  
50 transfer. We identified various genetic *pks* configurations indicative of an introduction of the  
51 *pks* island into *E. coli* on multiple independent occasions. Despite the existence of various  
52 acquisition scenarios, we found that the *pks* sequences were highly conserved and *pks*-  
53 carrying strains were overwhelmingly capable of producing colibactin, suggesting that the *pks*  
54 island is under selective pressure, through the production of colibactin or other secondary  
55 metabolites. Future implications include the identification of such metabolites and their  
56 biological activities that could be advantageous to *E. coli* and enable its adaptation to various  
57 ecological niches.

## 58 **DATA SUMMARY**

59 All sequence data of the 785 *E. coli* used in this study are freely available from the NCBI  
60 BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under the accession number  
61 PRJDB5579. This database was updated to include the sequence data obtained using ONT  
62 MinION for the *E. coli* reference strain SP15 and for *E. coli* strains ECSC054, JML285, KS-  
63 NP019, NS-NP030 and SI-NP020. The sequence data of *E. coli* strain UPEC129 obtained using  
64 PacBio instrument were deposited in the NCBI BioProject database and are available at  
65 <https://www.ncbi.nlm.nih.gov/Traces/study/> under the accession number PRJNA669570.  
66 Hybrid MinION-Illumina and PacBio-Illumina assemblies are available at the NCBI nucleotide

67 database. The genome sequences of 36 other *E. coli* reference strains and 7 non-*E. coli* strains  
68 were retrieved from NCBI.

## 69 INTRODUCTION

70 *Escherichia coli* is not only a commensal resident of the human and animal gut, but also a  
71 pathogen responsible for intestinal or extra-intestinal infections. The *E. coli* species is  
72 characterized by a high genetic and phenotypic diversity, with a population distributed into at  
73 least eight major phylogenetic groups (A, B1, B2, C, D, E, F and G) (1). *E. coli* strains from  
74 the phylogroup B2 are increasingly found in the feces of healthy humans in high-income  
75 countries and also responsible for extra-intestinal diseases, including urinary tract infections,  
76 sepsis, pneumonia and neonatal meningitis (2). By enabling the exchange of genetic material  
77 between bacterial cells, horizontal gene transfer (HGT) is a major driving force in the evolution  
78 of bacteria, including adaptation to their host and expansion of their ecological niche (3). HGT-  
79 mediated acquisition of large genomic islands (GIs) or pathogenicity islands (PAIs) is  
80 recognized as a major contributor to the emergence of the various *E. coli* pathotypes (4). The  
81 *E. coli pks* pathogenicity island consists of a *clbA-clbS* gene cluster enabling the biosynthesis  
82 of a polyketide (PK) - non-ribosomal peptide (NRP) hybrid genotoxin known as colibactin (5).  
83 This island exhibits typical features of horizontally acquired genomic elements: (i) it is a large  
84 (*i.e.* 54-kb) region with a distinct GC content compared to that of the chromosomal backbone,  
85 (ii) it is physically associated with a phage-type integrase gene that probably mediated its  
86 insertion into the chromosome, and (iii) it is located at a tRNA *locus* and is flanked by two short  
87 (*i.e.* 17-bp) direct repeats (DRs) reminiscent of those generated upon integrase-mediated  
88 insertion of mobile genetic elements (5, 6). The *pks* island can be found in other members of  
89 *Enterobacteriaceae* such as *Klebsiella pneumoniae*, *Citrobacter koseri* and *Enterobacter*  
90 *aerogenes* (6), and in the honeybee gut commensal *Frischella perrara* (7) and the marine  
91 sponge commensal *Pseudovibrio* sp. (8).  
92 Colibactin is a virulence factor for extra-intestinal pathogenic *E. coli* (ExPEC) (9-11) and is  
93 also a suspected procarcinogenic factor (12-14). Colibactin induces DNA interstrand cross-

94 links (ICLs) (15) and double-strand breaks (5) in host eukaryotic cells. Its production involves  
95 the sequential action of the Clb proteins, including PK synthases (PKSs), NRP synthetases  
96 (NRPSs), hybrid PKS-NRPS and accessory, editing and maturation enzymes (16). Colibactin  
97 is first synthesized as a prodrug called precolibactin, carrying an N-myristoyl-D-Asparagine  
98 (C14-Asn) side chain that is then cleaved in the periplasm to release the active genotoxin, whose  
99 translocation across the bacterial outer membrane remains unknown (17). The production of  
100 colibactin is positively regulated by ClbR (18). The multi-modular PKS-NRPS assembly line  
101 not only produces colibactin but also a set of numerous secondary metabolites with varying  
102 modes of action (19, 20). These include analgesic lipopeptides, such as C12-Asn-GABA, with  
103 the capability to diffuse across the epithelial barrier and act on sensory neurons to decrease  
104 visceral pain in the host (21). The *pks* island also contributes to the production of siderophores  
105 (enterobactin, salmochelin and yersiniabactin), *via* its promiscuous phosphopantetheinyl  
106 transferase ClbA (10), and siderophore-microcins *via* its ClbP peptidase (22).  
107 To date, the presence of the *pks* island was investigated mostly in *E. coli* strains isolated from  
108 humans with extra-intestinal infections (5, 6, 23, 24). Here we explored the distribution,  
109 conservation and functionality of the *pks* island in a large collection of non-clinical *E. coli*  
110 strains originating from human and bovine hosts (25). We found that the *pks* island was not  
111 only present in the phylogroup B2 but also in other genetic phylogroups. We identified different  
112 scenarios for its integration into the *E. coli* genome. The sequence of the *pks* island is highly  
113 conserved and *pks*-positive strains were overwhelmingly capable of producing colibactin,  
114 suggesting that the *pks* island is under selective pressure for the adaptation of *E. coli* to various  
115 ecological niches, through the production of colibactin or other metabolites or *pks*-encoded  
116 enzymatic activities.

117

118

## 119 **METHODS**

### 120 **Bacterial strains used in the study**

121 The *E. coli* strains were collected in Japan from 418 healthy bovines in 2013 and 2014, 278  
122 healthy humans in 2008, 2009 and 2015, and 89 humans with extra-intestinal infections, either  
123 bacteremia (n=67) in 2002-2008 or urinary tract infection (n=22) in 2006 and 2011. They were  
124 described recently (25) and corresponded each to a single isolate, duplicates showing less than  
125 5 SNPs difference in their whole genomes being excluded from this study. A list of the 109 *pks*-  
126 positive isolates is provided in Table S1. Additional 37 *E. coli* reference strains (Table S2) and  
127 7 non-*E. coli* strains (Table S3) were included in this study; their genome sequences were  
128 downloaded from NCBI, except for *E. coli* SP15 which was not available and was obtained  
129 here (see below).

### 130 **Whole genome sequencing**

131 The whole genome sequences of the 785 *E. coli* isolates were determined by Illumina  
132 sequencers (25). Among these, the genomes of 5 *E. coli* strains (ECSC054, JML285, KS-  
133 NP019, NS-NP030, and SI-NP020) were further subjected here to long-read sequencing using  
134 Oxford Nanopore Technologies (ONT) MinION device. The DNA libraries were prepared  
135 using the rapid barcoding kit (Oxford Nanopore Technologies) and sequenced  
136 using MinION R9.4.1 flow cells. Long-read sequencing of *E. coli* strain UPEC129 was also  
137 performed using Pacific Biosciences (PacBio) RSII sequencer (Genoscreen, Lille, France). The  
138 DNA was extracted using Genra Puregen Yeast/Bact (Qiagen) and the DNA libraries prepared  
139 using the SMRTbell Template Prep kit (PacBio). Hybrid assembly of Illumina paired-end reads  
140 and MinION or PacBio reads was performed using Unicycler (v.0.4.8) (26). The whole genome  
141 sequence of *E. coli* reference strain SP15 was obtained using Illumina and ONT MinION  
142 instruments and assembled as described above.

143 **Sequence and phylogenetic analysis**

144 The core gene-based phylogenetic tree was constructed as described previously (25). Briefly,  
145 core genes were determined using Roary (27) and SNP sites were extracted from the core gene  
146 alignment using SNP-sites (28). The maximum likelihood (ML) tree was constructed using  
147 RAxML (29) with the GTR-GAMMA model and displayed using iTOL (30).

148 For the phylogenetic analysis of the entire *pks* island, the genome sequences of *pks*-positive  
149 strains were aligned with the entire *pks* island sequence of strain IHE3034 using MUMmer (31)  
150 and the SNP sites located therein were identified. After removing SNP sites on the VNTR  
151 region, a neighbor-joining (NJ) tree was constructed by MEGA7 (32) using the Tamura-Nei  
152 evolutionary model.

153 Cophylogenetic analysis of the core-gene based ML tree and the *pks*-based NJ tree was  
154 performed using the “cophylo” function of the R package Phytools (33).

155 Sequence type and phylogroup determination was performed as described previously (25).

156 The *pks* sequences from four *E. coli* strains belonging to distinct phylogroups (i.e. SI-NP020,  
157 KS-NP019, UPEC129 and ECSC054 from phylogroups A, B1, B2 and D, respectively) were  
158 extracted from hybrid assemblies and compared at the nucleotide level with that of the reference  
159 *E. coli* strain IHE3034. In addition, the amino acid sequences were obtained for the 19 *clb* genes  
160 of each strain and aligned by MUSCLE with MEGA7 (32). The alignment file was analyzed  
161 with the sequence identity and similarity online software  
162 (<http://imed.med.ucm.es/Tools/sias.html>; accessed in July 2020).

163 The comparison of *pks* sequences from *E. coli* and other bacterial species was performed with  
164 BLASTn ([https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE\\_TYPE=BlastSearch](https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastSearch)). Each *pks*  
165 region was defined from *clbA* to *clbS* and used as the query nucleotide sequence against each  
166 *pks* region as the subject. Then, the alignment was visualized with Artemis Comparison Tool  
167 (v13.0.0) (34).



168 The integrase nucleotide and amino-acid sequences were aligned using MUSCLE (v3.8.31) and  
169 the phylogeny was analysed with PhyML (v3.1/3.0 aLRT) prior to tree visualisation with  
170 TreeDyn (v198.3) (<http://www.phylogeny.fr>; accessed in Sept 2020).

171 The CC95 strains were typed for their *fimH* allele using FimTyper (v1.0)  
172 (<https://cge.cbs.dtu.dk/services/FimTyper/>; accessed in Nov 2020) and were further assigned to  
173 subgroups A-E by analysis of the presence of either of the five subgroup-specific genes  
174 described previously (35).

### 175 **PCR analysis of the *clbJK* fusion gene**

176 The 5,651-bp deletion in the *clbJ-clbK* region resulting in the *clbJK* fusion gene was tested  
177 using a duplex PCR assay, with two primer pairs. The first primer pair (*clbK-F*, 5'-  
178 GACTGCCCAACATACGCTCCG-3'; *clbK-R*, 5'-TTGTGTCGTTGTACTCTCGGC-3') was  
179 used to amplify a 722 bp-long DNA fragment that is located within the deleted region and is  
180 thus only present in strains with an intact *clbJ-clbK* region. The second primer pair consisted  
181 of primers *clbJK-F* (5'-AGAATTACCCACTGCCACCA-3') and *clbJK-R* (5'-  
182 GGCGCTAATGGATCAGATGT-3') flanking the deleted region, and was used to amplify a  
183 1441 bp-long DNA fragment only present in strains with a *clbJK* fusion gene. The strains with  
184 an intact *clbJ-clbK* region or a *clbJK* fusion gene yielded a 722-bp or a 1441-bp long  
185 amplification product, respectively. Reaction mixture of 50  $\mu$ L final volume contained 2 $\mu$ L  
186 template DNA, 1X GoTaq Reaction buffer, 200 $\mu$ M of each dNTP, 4 mM of MgCl<sub>2</sub>, 1.25 U of  
187 GoTaq DNA polymerase (Promega, France) and 0.2 $\mu$ M of each primer (Eurofins Genomics  
188 Ebersberg, Germany). Amplification was done in a GeneAmp® 9700 thermal cycler (Applied  
189 Biosystems, Courtaboeuf, France), with the following program: initial denaturation at 95°C for  
190 2 min; 30 cycles of denaturation at 95°C for 30 s, annealing at 56°C for 45 s and extension at  
191 72°C for 1 min 30 s; final extension at 72°C for 5 min. Electrophoresis was carried out in 1%  
192 agarose gel and the PCR products visualized after Gel Red (Biotium) staining using a Bio-Rad

193 Chemidoc XRS system (Bio-Rad, France).

#### 194 ***In vitro* DNA interstrand crosslinking assay**

195 ICL activity was assessed as described previously (15). Briefly, 3 10<sup>6</sup> *E. coli* cells or 6 10<sup>6</sup>  
196 *Erwinia oleae* cells pre-grown for 3.5 h in DMEM with 25 mM HEPES (Invitrogen) were  
197 mixed with EDTA (1 mM) and 400 ng of linearized plasmid pUC19 DNA and the mixtures  
198 were incubated for 40 min at 37°C. After pelleting the bacteria, the DNA was purified from the  
199 supernatant and analyzed by electrophoresis on denaturing (40 mM NaOH - 1 mM EDTA) 1%  
200 agarose gels. ICL activity of *E. oleae* was also tested in the presence of 400nM 6-histidine-  
201 ClbS, which was purified with HisPur nickel-nitrilotriacetic acid (Ni-NTA) agarose (Thermo  
202 Scientific) from a culture of BL21(DE3) strain hosting the plasmid pET28a-ClbS-His, as  
203 described previously (15).

#### 204 **Megalocytosis assay**

205 Non-hemolytic *pks*-positive strains were tested for megalocytosis on infected HeLa cells as  
206 described previously (5, 36). Briefly, HeLa cells grown to 50% confluence in cell culture 96-  
207 well plates were inoculated with 5 µL of overnight culture of bacteria in infection medium  
208 (DMEM with 25 mM HEPES) and incubated for 4 h at 37°C in a 5% CO<sub>2</sub> atmosphere. Cells  
209 were then washed and incubated 48 to 72 h in cell culture medium supplemented with 200  
210 µg/mL gentamicin, and then stained with methylene blue for microscopy examination.

#### 211 **H2AX phosphorylation assay**

212 HeLa cells were infected as described above and H2AX phosphorylation was quantified  
213 immediately after the 4 h infection step by immunofluorescence as described elsewhere (37).

#### 214 **C14-asn quantification**

215 *E. coli* strains were grown for 24 h at 37°C in 10 mL DMEM-HEPES (Gibco), resuspended in  
216 500µL HBSS (Invitrogen) and then crushed with a Precellys instrument (Ozyme, Montigny le

217 Bretonneux, France). After addition of an internal standard mixture (Deuterium-labeled  
218 compounds; 400 ng/mL), cold methanol (MeOH) was added and samples were solid-phase  
219 extracted on HLB plates (OASIS® HLB 2 mg, 96-well plate, Waters, Ireland). Lipids were  
220 eluted with MeOH, evaporated under N<sub>2</sub>, resuspended in MeOH and analysed by high-  
221 performance liquid chromatography/tandem mass spectrometry analysis (LC-MS/MS)  
222 (MetaToulLipidomics Facility, INSERM UMR1048, Toulouse, France), as described  
223 previously (21).

## 224 RESULTS

### 225 *The pks island was mainly found in specific E. coli lineages from phylogroup B2*

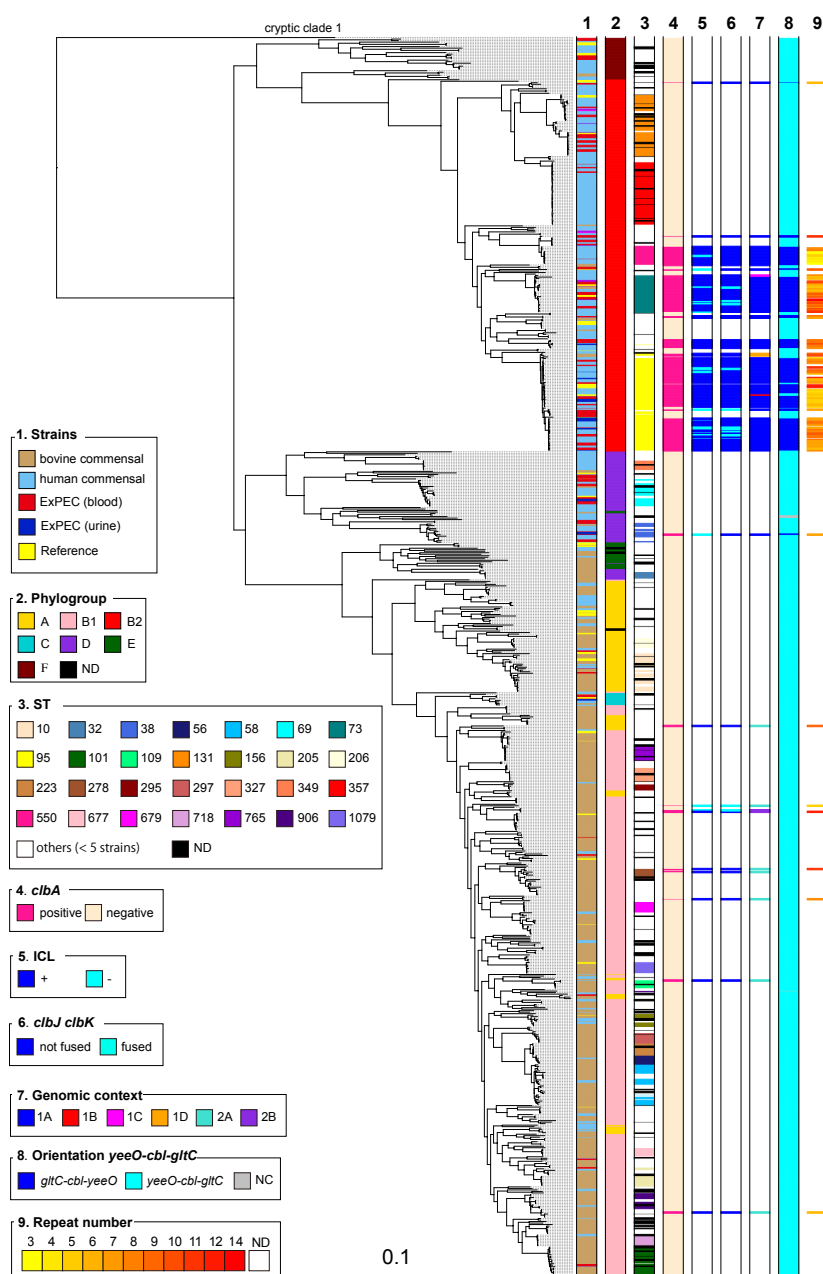
226 The presence of the *pks* island was investigated in a collection of 785 *E. coli* strains (25)  
227 belonging to at least 296 different sequence types (STs) and originating mostly from fecal  
228 samples of healthy bovines and humans. Clinical isolates recovered from urine or blood  
229 samples of human patients with extra-intestinal infection were also included for comparison.  
230 We detected the *pks* island in 109 *E. coli* strains, including 62 (22.3%) out of 278 healthy human  
231 fecal isolates and 12 (2.9%) out of 418 healthy bovine fecal isolates (Table 1; Fig. 1). As  
232 expected a higher proportion of *pks*-positive strains were found among ExPEC, i.e. 35 (39%)  
233 out of 89 strains, including 14 (63.6%) out of 22 strains from urinary tract infection and 21  
234 (31.3%) out of 67 strains from bacteremia. The vast majority of the 109 *pks*-positive strains  
235 corresponded to B2 isolates (Table 1; Fig. 1) and the *pks* island was mainly present in specific  
236 lineages or STs of the B2 phylogroup (Fig. 1).

237 **Table 1.** Occurrence of *pks* in *E. coli* strains from healthy humans or bovines, and human  
238 patients with extra-intestinal infection.

240 Origin	241 Phylogroup (no. <i>pks</i> + strains / no. strains tested)								
	242 A	B1	B2	C	D	E	F	Uncl <sup>c</sup>	Total
244 Healthy bovines <sup>a</sup>	1/48	7/314	4/11	0/20	0/12	0/8	0/4	0/1	12/418
246 Healthy humans <sup>a</sup>	0/14	1/29	61/163	0/11	0/37	0/1	0/20	0/3	62/278

248	Human patients <sup>b</sup>	0/3	0/5	34/59	0/4	1/10	0/0	0/8	0/0	35/89
249										
250	Total	1/65	8/348	99/233	0/35	1/59	0/9	0/32	0/4	109/785
251										

252 <sup>a</sup> Isolates were collected from feces of healthy individuals  
 253 <sup>b</sup> Isolates were collected from blood (n=67) or urine (n=22) from human patients with extra-intestinal infection.  
 254 <sup>c</sup> Unclassified



255

256

257 **Figure 1.** Phylogenetic relationship and distribution of *pks*-positive/negative *E. coli* isolates  
258 among 696 human and bovine commensal *E. coli*, 89 ExPEC and 37 completely sequenced  
259 reference *E. coli* strains. A core gene-based maximum likelihood (ML) tree was constructed  
260 based on 271,403 SNPs located on 2,000 core genes and rooted on cryptic *Escherichia* clade 1  
261 strains as outgroups. Origin (column 1), phylogroup (column 2), major sequence type (ST) (*i.e.*  
262 ST identified for at least 5 strains) (column 3), presence of *pks* (*clbA*) (column 4), colibactin  
263 activity (ICL) (column 5), presence of the *clbJK* fusion gene (column 6), genetic *pks*  
264 configuration (see Fig. 3) (column 7), orientation of the *asnV-asnU-asnW* region situated  
265 downstream the *asnT* tRNA gene (column 8) and the number of repeats 5'-ACAGATAC-3'  
266 found in the *clbB-clbR* intergenic region (see Fig. 2) (column 9) are shown for each strain. ND,  
267 not determined.

269 Strikingly, the *pks* island was found in (nearly) 100% of strains belonging to ST12, ST73, ST95  
270 and ST550, while it was excluded from other STs, such as ST131 and ST357 (Fig. 1; Table 2).

271 Interestingly, these *pks*-positive and -negative STs are found in distinct clusters in the core-  
272 genome based phylogenetic tree (Fig. 1) suggesting that *pks* acquisition occurred after the  
273 divergence of these clusters from a common ancestor. We further characterized the 54 *pks*-  
274 positive strains of ST95 for their *fimH* allele and affiliation to CC95 subgroups A to E defined  
275 previously (35). We could assign 35 of them to subgroup A (n=22), B (n=12) or E (n=1) (Table  
276 S1). The remaining 19 strains, including 15 of serotype O1:H1, did not belong to any of these  
277 five subgroups. No *pks*-positive strain was assigned to CC95 subgroups C or D, in agreement  
278 with previous results (35).

279 Except for four B2 strains originating from healthy bovines, the *pks*-positive B2 isolates  
280 originated from humans, either patients with extra-intestinal infection (n=34) or healthy  
281 individuals (n=61) (Fig. 1; Table 1). The low occurrence of *pks* among bovine isolates likely  
282 reflected the low prevalence of B2 strains in cattle (25). Interestingly, 10 non-B2 *pks*-positive  
283 strains were identified corresponding to 1 human blood isolate from phylogroup D, 1 healthy  
284 human fecal isolate from group B1 and 8 healthy bovine fecal isolates from groups A (n=1) and  
285 B1 (n=7) (Fig. 1, Table 1). In contrast to the B2 *pks*-positive isolates, these strains were  
286 scattered throughout the core genome phylogenetic tree and were not representative of any  
287 particular lineage or ST (Fig. 1; Table 2).

288 **Table 2.** Distribution of the *pks* island in the predominant sequence types (ST) among *E. coli*  
 289 strains isolated from healthy bovines, healthy humans and human patients with extra-intestinal  
 290 infection.

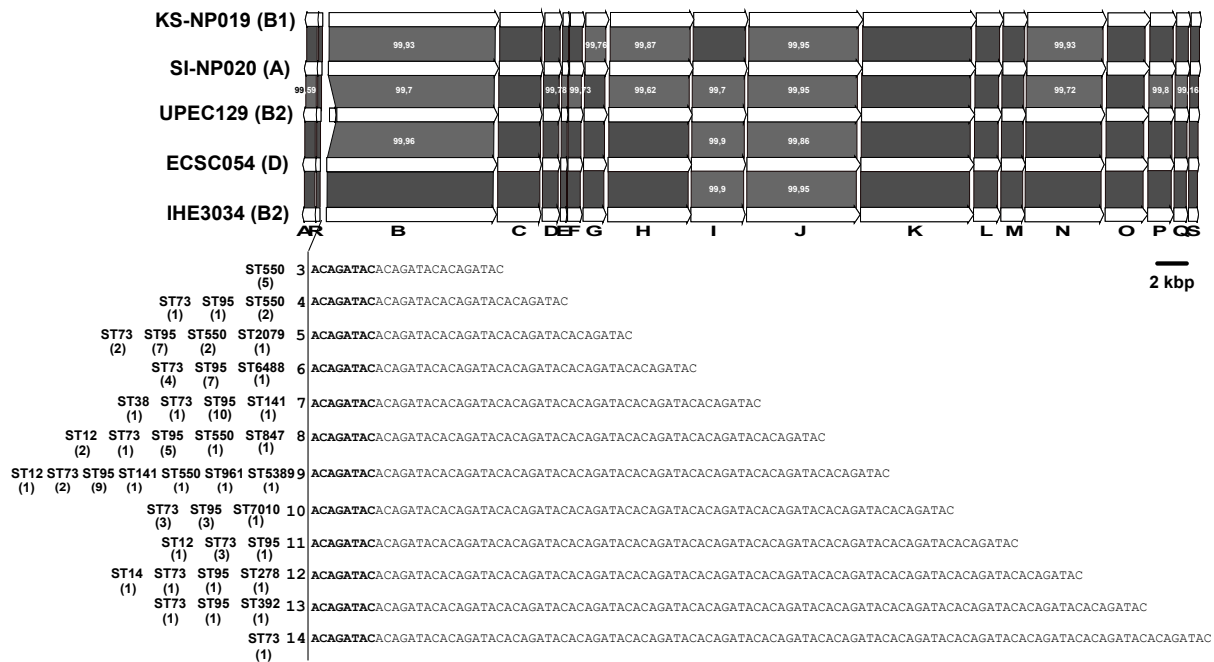
291	292	293	294
Phylogroup	ST <sup>a</sup>	<i>pks</i> <sup>+</sup> / no. strains	
294	A	10	0/22
295		206	0/5
296		6126	0/4
297			
298	B1	20	0/5
299		29	0/4
300		56	0/6
301		58	0/23
302		101	0/16
303		109	1/6
304		154	0/6
305		155	0/4
306		156	0/4
307		164	0/6
308		205	0/9
309		223	0/5
310		278	2/8
311		295	0/4
312		297	0/7
313		300	0/4
314		327	0/9
315		332	0/4
316		446	0/4
317		677	0/6
318		679	0/7
319		718	0/6
320		765	0/9
321		795	0/4
322		906	0/12
323		1079	0/7
324		1423	0/5
325		5487	0/4
326			
327	B2	12	4/4
328		73	20/20
329		95	54/56
330		131	0/25
331		357	0/35
332		550	12/12
333		1193	0/4
334			
335	C	88	0/4
336			
337	D	32	0/5
338		38	1/8
339		69	0/14
340		349	0/6

341  
 342 <sup>a</sup> only ST including at least 4 strains are listed  
 343

345 ***High level of genetic conservation of the pks island among E. coli phylogroups and other***  
346 ***enterobacteria***

347 The *pks* sequence from the B2 reference *E. coli* strain IHE3034 was compared to that of three  
348 non-B2 *E. coli* isolates, including the single group A isolate (i.e. SI-NP020), the single group  
349 D isolate (ECSC054) and one out of the eight B1 isolates (i.e. KS-NP019). An additional B2  
350 isolate (i.e. UPEC129) was also selected for this analysis. To perform this comparison, the  
351 whole genomes of these four isolates were assembled from a combination of short and long  
352 reads. At the amino acid level, over 99 % identity was observed for each of the 19 *clb* gene  
353 products (Fig. 2). At the nucleotide level, the only variation observed in the *pks* sequence was  
354 the size of the region located between *clbB* and *clbR* which contains a variable number of  
355 tandem repeats (VNTR) of the motif 5'-ACAGATAC-3' (6). This VNTR locus contained  
356 between 3 and 14 repeat units when the whole collection of *pks*-positive strains was analysed  
357 (except for 17 isolates for which the VNTR length could not be calculated), with no apparent  
358 correlation with the STs (Fig. 2). Therefore, apart from the size of the VNTR, the *pks* island  
359 was highly conserved among the strains, irrespective of their phylogroup or ST.

360



361

362 **Figure 2.** Comparison of the *pks* islands of *E. coli* strains belonging to phylogroups A, B1, B2  
 363 and D. The 19 ORFs of the *clbA-clbS* gene cluster from the reference *E. coli* strain IHE3034  
 364 sequence (group B2) and MinION- or PacBio-derived sequences of *E. coli* strains KS-NP019  
 365 (group B1), SI-NP020 (group A), UPEC129 (group B2) and ECSC054 (group D) are  
 366 represented by arrows with the arrowhead representing direction of transcription. The areas  
 367 between the corresponding genetic maps shaded in dark and light gray indicate 100% amino  
 368 acid identity and *ca.* 99% amino acid similarity, respectively. The number of repeated motif 5'-  
 369 ACAGATAC-3' found in the *clbB-clbR* nucleotide intergenic region of *pks*-positive *E. coli*  
 370 strains and the sequence types (ST) of the corresponding strains are indicated below the *clbA-*  
 371 *clbS* gene cluster, with the number of strains into parenthesis.  
 372

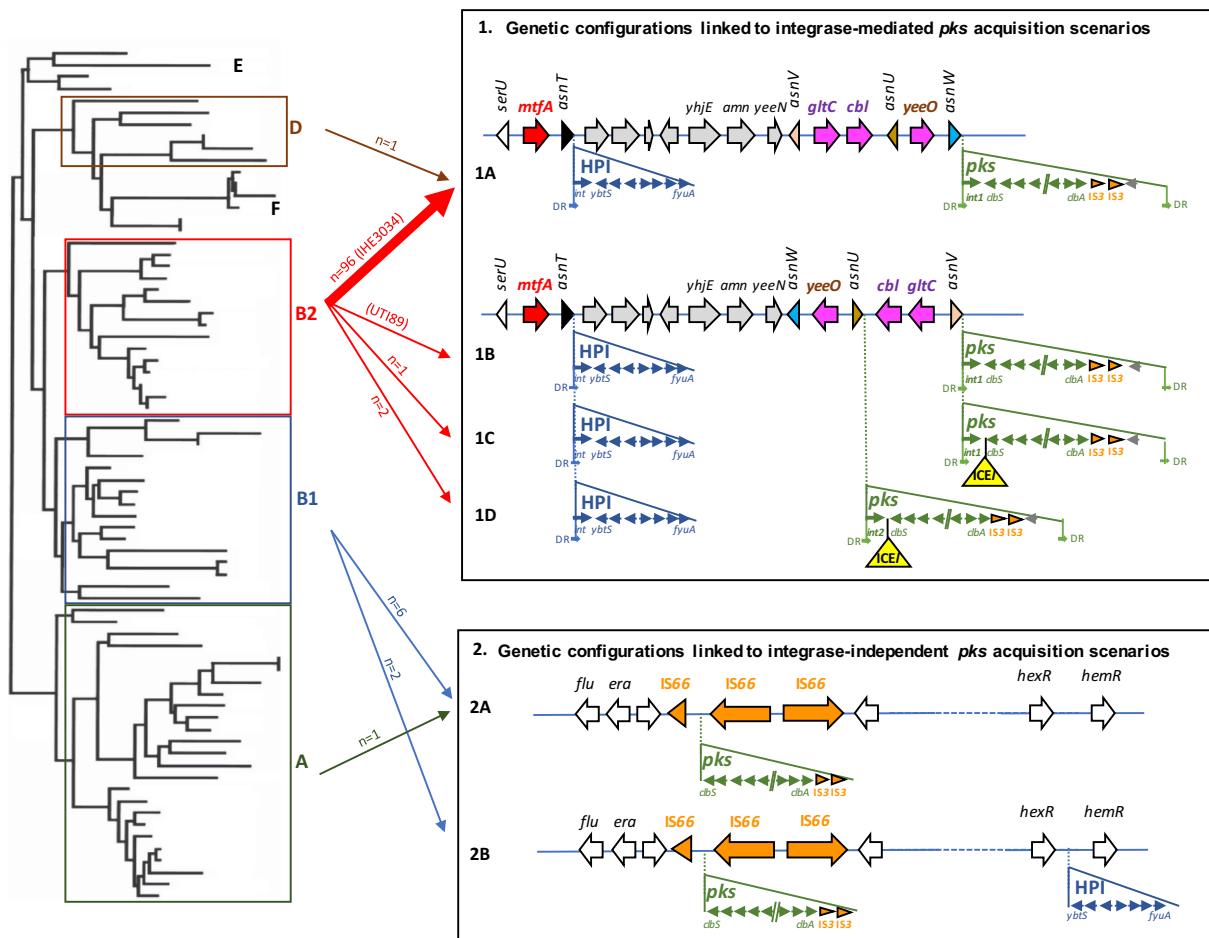
373 Comparison of the *pks* island nucleotide sequence from B2 reference strain IHE3034 with that  
 374 of other *pks*-positive bacterial species confirmed that it was conserved in other members of the  
 375 *Enterobacteriaceae* (Fig. S1) such as *K. pneumoniae*, *E. aerogenes*, *C. koseri*, *Serratia*  
 376 *marcescens* and *Erwinia oleae*. A similar *pks* island was present, although less conserved in *F.*  
 377 *perrara* and *Pseudovibrio* sp. (Fig. S1).

378



379 *The pks islands in E. coli from phylogenetic groups B2 and D share a similar genomic*  
 380 *environment.*

381 To gain insights into the events leading to the acquisition of the *pks* island into the *E. coli*  
 382 population, we analysed the genomic environment of the *pks* island in the 109 *pks*-positive  
 383 strains and in seven *pks*-positive *E. coli* reference strains (536, ABU83972, CFT073, Nissle  
 384 1917, UTI89, IHE3034 and SP15). Various configurations were found for the *pks* island  
 385 genomic environment, suggesting two main scenarios of *pks* acquisition, depending on the  
 386 presence or absence of an integrase gene (Fig. 3). The genetic configuration typical of B2  
 387 strains, named 1A, which is characterized by a *pks* island carrying an integrase gene and  
 388 inserted into the *asnW* tRNA gene in the vicinity of the *asnT*-located High Pathogenicity Island  
 389 (HPI) (5, 6) was found in 96 B2 strains of our collection and in one phylogroup D strain,  
 390 ECSC054 (Fig 1; Fig. 3).



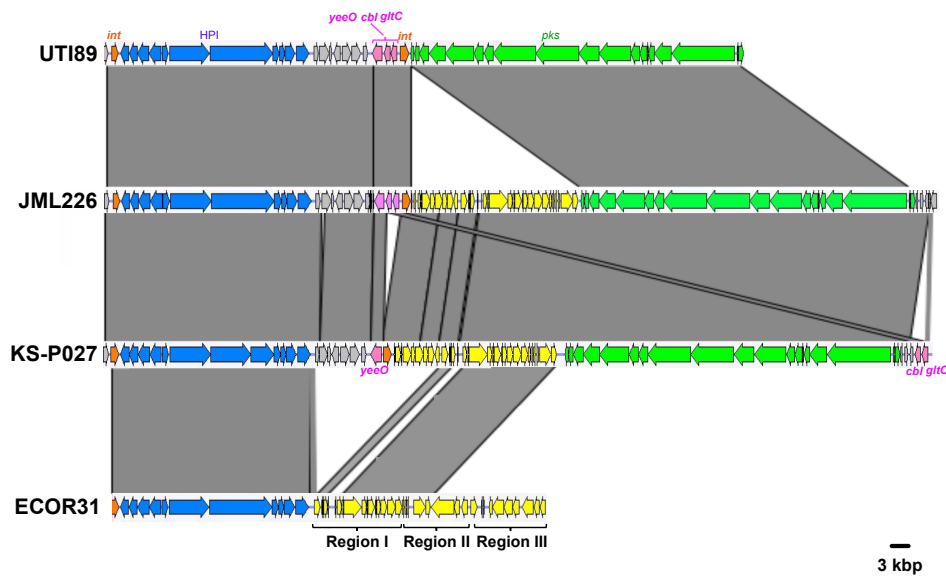
391

392 **Figure 3.** Genetic configurations of the *pks* island and HPI in *E. coli* strains and proposed  
393 scenarios for their acquisition. **Left.** Schematic phylogenetic tree showing the distribution of *E.*  
394 *coli* into the main phylogroups. **Right.** *E. coli* genetic *pks* and HPI configurations resulting from  
395 proposed acquisition scenarios involving site-specific recombination (configurations 1A, 1B,  
396 1C and 1D) or not (configurations 2A and 2B). The location and orientation of the tRNA genes  
397 and the ORFs of the chromosomal regions, including the integrase and genes from *pks* and the  
398 HPI, are indicated by the arrows. Partial and complete IS elements are represented by orange  
399 arrowheads and arrows, respectively. The ICE-like element (*ICEI*) found in configurations 1C  
400 and 1D is represented as a yellow triangle. DR, direct repeats located at the extremities of the  
401 islands (except for the HPI in configurations 1A-1D, one DR lacking at the right border).  
402 **Middle.** The arrows connect the phylogroups A, B1, B2 and D (left) with the *pks* and HPI  
403 configurations (right). The number of *E. coli* isolates belonging to the collection of 785 strains  
404 and corresponding to each configuration are indicated (except for configuration 1B which was  
405 only found in reference strain UTI89 indicated in parenthesis). The thick arrow represents the  
406 most frequently found configuration (exemplified here by reference strain IHE3034 indicated  
407 in parenthesis).

408

409 A similar configuration was found in a few other B2 strains but with variations in the location  
410 of the *pks* island which was inserted either into the *asnV* (corresponding to configurations 1B  
411 and 1C found in ST95 reference strain UTI89 and in ST73 strain JML226, respectively) or  
412 *asnU* tRNA gene (corresponding to configuration 1D found in strains KS-P003 and KS-P027,  
413 both belonging to ST95) (Fig. 1; Fig. 3). Besides the difference in the tRNA insertion site, the  
414 configuration 1B found in reference strain UTI89 differed from the major configuration 1A by  
415 the orientation of the 4,309-bp *asnW-asnU-asnV* tRNA region upstream of the *pks* island. This  
416 region contains three other genes, namely *gltC* and *cbl*, encoding two LysR-family  
417 transcriptional regulators, and *yeeO*, encoding a flavin mononucleotide [FMN] and flavin  
418 adenine dinucleotide [FAD] exporter. The configurations 1C and 1D possessed the same *asnW-*  
419 *asnU-asnV* orientation as in UTI89 and carried a 25-kb region between the *pks* integrase gene  
420 and *clbS*. A 14-kb section from this region exhibited high sequence similarity (>99%) to  
421 integrative and conjugative elements (ICE) identified in *E. coli* (*ICEEc1*) and *K. pneumoniae*  
422 (*ICEKp1*), in particular to the DNA regions I and II from *ICEEc1* involved in mating-pair  
423 formation (Mpf) and DNA mobilization, respectively (Fig. 4) (6, 38). This 14-kb section could

424 therefore be considered as an ICE-like element, although it is most likely non-functional given  
 425 the lack of a complete region II (Fig. 4). The remaining 11-kb section was not homologous to  
 426 ICE*Ec1*, ICE*Kp1* or any other ICE, and its role could not be predicted.



427  
 428 **Figure 4.** Comparison of the chromosomal region covering the HPI and *pks* island between 3  
 429 atypical *E. coli* B2 strains (UTI89, JML226 and KS-P027, with genetic configurations 1B, 1C  
 430 and 1D, respectively) and the integrative conjugative element ICE*Ec1* from *E. coli* strain  
 431 ECOR31. Nucleotide sequence similarity (>99%) between different DNA regions is indicated  
 432 by grey areas between the corresponding genetic maps. The *pks* island and the HPI are  
 433 represented in green and blue, respectively, and the integrase genes in orange. The *yeeO*, *cbl*  
 434 and *gltC* genes located between the *asnV* and *asnW* tRNA genes in UTI89 are represented in  
 435 pink. The region between the HPI and the *yeeO* gene is represented in grey and the ICE-related  
 436 region inserted either next to *pks* (JML226 and KS-P027) or next to the HPI (ECOR31) is  
 437 represented in yellow. In the ECOR31 strain, the ICE is divided in three parts, including region

438 I encoding a mating pair formation system, region II encoding a DNA-processing system, both  
439 involved in conjugative transfer, and region III comprising hypothetical genes.

440 The phage-type *pks* integrase is a tyrosine site-specific recombinase with similarity to the phage  
441 P4 integrase C-terminal catalytic domain (INT\_P4\_C). The integrase genes located at the *asnW*  
442 (configuration 1A) or *asnV* loci (configurations 1B and 1C) and their gene products were highly  
443 conserved and grouped into the integrase family 1 (Fig. S2), whereas the integrase genes located  
444 at the *asnU* locus (configuration 1D) and their gene products shared 94% nucleotide and 94.6%  
445 amino acid sequence similarity, respectively, with those of family 1 and were thus grouped into  
446 the integrase family 2 (Fig. S2).

447

#### 448 ***Atypical genomic environments of pks islands in E. coli from phylogenetic groups A and B1***

449 In the 9 *pks*-positive *E. coli* strains from phylogroups A and B1, two different configurations  
450 (named 2A and 2B) were observed that drastically differed from those found in B2/D strains.  
451 Their *pks* islands lacked an integrase gene, were not inserted into a tRNA gene and there were  
452 no direct repeats at their chromosomal boundaries (Fig. 3). The *pks* islands were located in the  
453 vicinity of the genes *flu* (or *agn43*) and *era* encoding the Ag43 autotransporter adhesin and a  
454 GTPase essential for cell growth and viability, respectively. They were flanked on one side by  
455 a truncated copy of the IS66 insertion sequence (IS) and on the other side by two intact IS66  
456 copies. Two truncated copies of IS3 were also found next to the *clbA* gene but this was also the  
457 case for configurations 1A-1D. Moreover, in these B1/A strains, the HPI was absent (Table 3;  
458 Fig. 3, configuration 2A), except for two isolates in which the HPI was present but not in the  
459 vicinity of the *pks* island and not into a tRNA locus (Table 3; Fig. 3, configuration 2B). Using  
460 PCR assays, it was shown previously that three *E. coli* strains from phylogroup B1 (namely  
461 U12633, U15156 and U19010) possessed a *pks* island that co-localized with the HPI and the

462 DNA transfer and mobilization region of an ICE*Ec1*-like element (6), a situation that is  
463 reminiscent of that of configuration 1D. However, as the whole genome sequences of these  
464 three strains were not available, this could not be confirmed here.

465 Since the *asnW-asnU-asnV* tRNA region displayed distinct orientations in *pks*-positive B2  
466 strains depending on *pks* configuration, we further analysed its orientation for the rest of the *E.*  
467 *coli* collection, i.e. in *pks*-positive B1/A strains and in *pks*-negative strains. The “*asnV-asnU-*  
468 *asnW*” orientation was uniquely found in typical *pks*-positive B2 strains with configuration 1A,  
469 suggesting that, in these strains, *pks* acquisition at the *asnW* locus was accompanied by an  
470 inversion of the upstream tRNA-encoding region (Fig. 1; Fig. 3).

471 **Table 3.** Characteristics of B2 and non-B2 *pks*-positive *E. coli* strains with atypical features regarding *pks* integrity, functionality or location.  
 472

473

474	Group	Strain	Origin	Sample	Year	ST	Serotype	<i>hlyA</i> (hemolysis)	<i>ybt</i> (locus)	<i>pks</i> (locus)	<i>clbJ-clbK</i> <sup>a</sup>	Megal. <sup>b</sup>	ICL	H2AX <sup>b</sup>	C14-Asn <sup>b,c</sup>
475	A	SI-NP020	b	feces	2014	7010	uncl:H14	- (-)	-	+ (not tRNA)	wt	+	+	+	+++
476	B1	JML285	H	feces	2015	109	Gp2:H8	- (-)	-	+ (not tRNA)	wt	+	+	+	+++
478		HH-NP008	b	feces	2014	847	uncl:H2	- (-)	-	+ (not tRNA)	wt	+	+	nt	nt
479		KK-NP025	b	feces	2014	6488	uncl:H8	- (-)	-	+ (not tRNA)	wt	nt	+	nt	nt
480		KS-NP019	b	feces	2013	392	O8:H2	+ (+)	+ (not tRNA)	+ (not tRNA)	wt	nt	+	nt	+++
481		SI-NP013	b	feces	2014	278	uncl:H7	- (-)	-	+ (not tRNA)	wt	+	+	nt	nt
482		SI-NP017	b	feces	2014	278	Gp2:H21	- (-)	-	+ (not tRNA)	wt	+	+	nt	nt
483		NS-NP014	b	feces	2014	2079	O8 :H19	- (-)	-	+ (not tRNA)	f	nt	-	nt	nt
484	B2	NS-NP030	b	feces	2014	392	uncl:H2	+ (+)	+ (not tRNA)	+ (not tRNA)	f	nt	-	nt	++
485		JML114	H	feces	2015	73	O6:H1	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	wt	-	-	-	++
486		JML165	H	feces	2015	550	uncl:H5	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	wt	-	-	nt	++
487		JML201	H	feces	2015	95	O1:H1	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	wt	-	-	nt	nt
488		JML226	H	feces	2015	73	Gp7:H12	+ (+)	+ ( <i>asnT</i> )	+ ( <i>asnV</i> )	wt	nt	+	nt	nt
489		KS-P003	b	feces	2013	95	Gp7:H5	+ (+)	+ ( <i>asnT</i> )	+ ( <i>asnU</i> )	wt	nt	+	nt	nt
490		KS-P027	b	feces	2013	95	Gp7:H5	+ (+)	+ ( <i>asnT</i> )	+ ( <i>asnU</i> )	wt	nt	+	nt	nt
491		JML008	H	feces	2015	95	Gp7:H4	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	nt	nt
492		JML102	H	feces	2015	73	O6:H1	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	nt	nt
493		JML282	H	feces	2015	95	Gp7:H7	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	nt	nt
494		JML288	H	feces	2015	95	O1:H7	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	nt	nt
495		JML291	H	feces	2015	95	O1:H12	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	nt	nt
496		JML296	H	feces	2015	73	uncl:H1	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	-	++
497		SI-NP032	b	feces	2014	73	O25:H5	+ (+)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	nt	-	nt	++
498		ECSC09	H	blood	2006	95	Gp7:H7	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	nt	nt
499		UPEC57	H	urine	2011	95	Gp7:H7	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	-	-	-	nt
500		UPEC91	H	urine	2011	95	O1:H7	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	f	nt	-	nt	nt
501	CM1	H	urine	2006	95	O1:H1	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	wt	-	-	nt	-	
502	UPEC129	H	urine	2011	uncl	Gp7:H7	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	wt	-	-	-	-	
503	D	ECSC054	H	blood	2004	38	O4:H30	- (-)	+ ( <i>asnT</i> )	+ ( <i>asnW</i> )	wt	-	-	-	+
504															

505

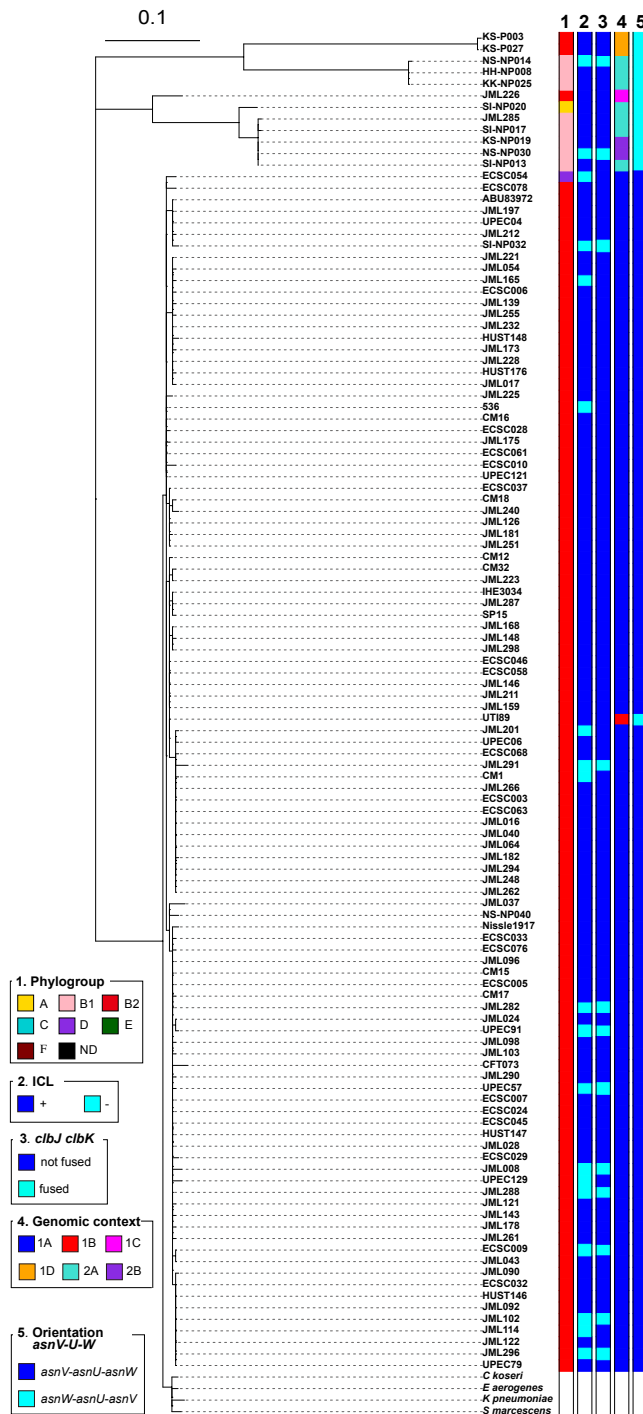
506  
 507 <sup>a</sup> wt, full length *clbJ* and *clbK* genes; f, *clbJK* fusion gene

508 <sup>b</sup> Megal., megalocytosis; nt, not tested

509 <sup>c</sup> -, no C14-Asn detected ; +, ca. 50-400 pg C14-Asn / 10e8 CFU; ++, ca. 400-600 pg C14-Asn / 10e8 CFU; +++, ca. 650-1200 pg C14-Asn / 10e8 CFU.

510 ***The phylogeny of the pks island globally reflects that of the E. coli core genome***

511 To shed further lights on the different *pks* acquisition scenarios, we constructed a phylogenetic  
512 tree of the entire *pks* sequences (i.e. from *clbA* to *clbS*, except for the VNTR-containing region  
513 which was excluded from the analysis) from the 109 *pks*-positive strains. Globally, the *pks*  
514 sequences from the strains showing distinct *pks* genomic configurations formed distinct clusters  
515 (Fig. 5). The *pks* sequence of strain UTI89 with a unique configuration (configuration 1B) was  
516 clustered with those of strains with configuration 1A (Fig. 5). Remarkably, the *pks* sequences  
517 of B1/A strains segregated separately from those of B2/D strains. Moreover, the *pks* sequences  
518 with an insertion of an ICE-like element in the B2 (ST73) human strain JML226 (with  
519 configuration 1C) or in the pair of B2 (ST95) bovine strains KS-P003 and KS-P027 (with  
520 configuration 1D) also clustered separately and were closer to the *pks* sequences of B1/A strains  
521 than to those of the B2/D strains lacking this ICE-like element (Fig. 5). Finally, the *pks*  
522 sequences of *C. koseri*, *E. aerogenes*, *K. pneumoniae* and *S. marcescens* were close to those of  
523 *E. coli* B2/D strains with configuration 1A (Fig. 5) whereas that of *E. oleae* was more  
524 phylogenetically distant and clustered separately (data not shown).



525

526 **Figure 5.** Phylogenetic tree of the entire *pks* island. SNP analysis was performed with the *pks*  
 527 sequences of the 109 *pks*-positive *E. coli* strains and a NJ phylogenetic tree was built. The *pks*  
 528 sequences from 7 reference *E. coli* strains (536, ABU83972, CFT073, Nissle 1917, UT189,  
 529 IHE3034 and SP15) and other enterobacteria (i.e. *C. koseri* ATCC BAA-895, *E. aerogenes*  
 530 EA1509E, *K. pneumoniae* 1084, and *S. marcescens* AS012490) were also included in this tree.  
 531 Phylogroup (column 1), colibactin activity (column 2), presence of a *clbJK* fusion gene (column  
 532 3), genetic *pks* configuration (see Fig. 3) (column 4) and orientation of the *asnV-asnU-asnW*  
 533 region (see Fig. 3) (column 5) are indicated for each *pks*-positive strain.

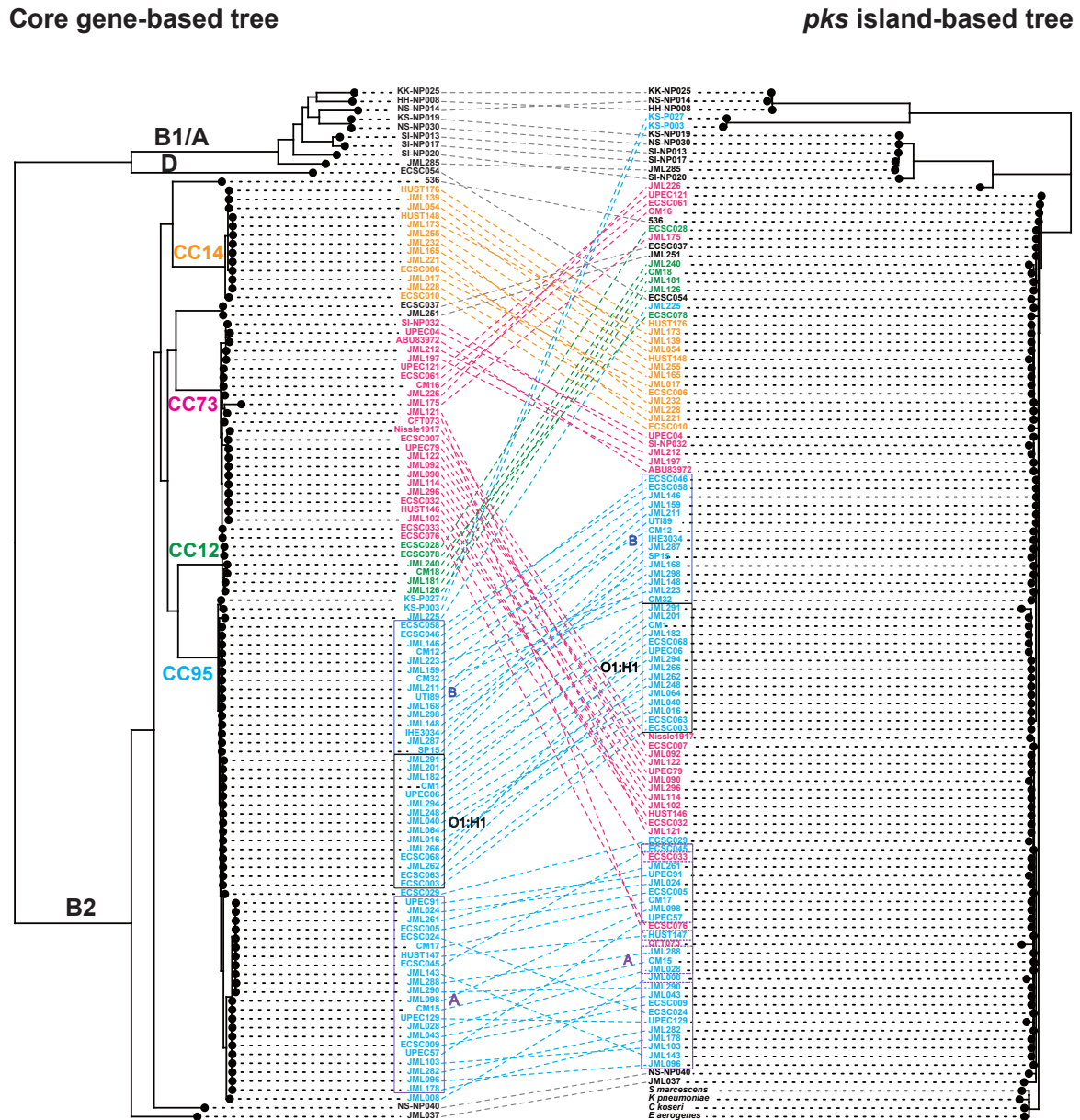
534



535 To further assess the evolutionary relationships between the *pks* sequences and the genetic  
536 background of the strains, a cophylogenetic analysis was performed where the phylogenetic  
537 trees based on the *pks* sequence and the core genome were compared. Globally, congruence  
538 was observed between both trees (Fig. 6). It was noticeable that most of the typical *pks*-positive  
539 B2 strains whose core genomes clustered together into lineages of clonal complexes (CC) 12,  
540 CC14, CC73 and CC95 contained *pks* sequences that also clustered together in different  
541 subgroups of the main *pks* cluster (Fig. 6). In particular, the CC95 strains that clustered together  
542 into subgroups A and B (as defined by *fimH* typing) or O1:H1 subgroup in the core genome  
543 tree also clustered together in the *pks* tree (Fig. 6). These observations support the hypothesis  
544 of an introduction of the *pks* island into CC12, CC14, CC73 and CC95 through horizontal  
545 acquisition by their most recent common ancestor (MRCA) or by the MRCA of each of these  
546 lineages, followed by vertical transmission with subtle *pks* divergence overtime. Since CC95  
547 subgroups C or D contain *pks*-negative strains (35) (strain ECSC026, subgroup C; this study),  
548 we further hypothesize that *pks* was lost during the evolution of these sublineages.

549 The fact that a single *pks*-positive strain from phylogroup D (ECSC054) possessed a *pks* island  
550 whose sequence clustered with that of B2 strains (Fig. 6) suggests that this strain acquired *pks*  
551 from a B2 strain through HGT. The *pks*-carrying B1/A strains were diverse based on their core  
552 genomes and their *pks* sequences clustered in two separate groups that were distantly related to  
553 the major *pks* cluster of B2 strains (Fig. 6), suggesting the existence of sporadic *pks* introduction  
554 within the B1 and A phylogroups, presumably through HGT from a donor strain different from  
555 typical *pks*-positive B2 strains. Finally, the cophylogeny also confirmed that the two atypical  
556 B2 ST95 strains KS-P03 and KS-P027 clustered with the other B2 strains of ST95 based on the  
557 core genome but contained a divergent *pks* sequence which was closer to those of B1 or A  
558 strains (Fig. 6), suggesting that this pair of strains probably acquired their *pks* islands through  
559 HGT, possibly from a donor strain carrying a *pks* island with an ICE insertion. The same

560 scenario also presumably occurred with the atypical B2 ST73 strain JML226 which clustered  
 561 with the other B2 strains of ST73 in the core genome tree but carried a *pks* island characterized  
 562 by an ICE-like insertion and a sequence closer to those of B1/A strains than to those of B2 ST73  
 563 strains.



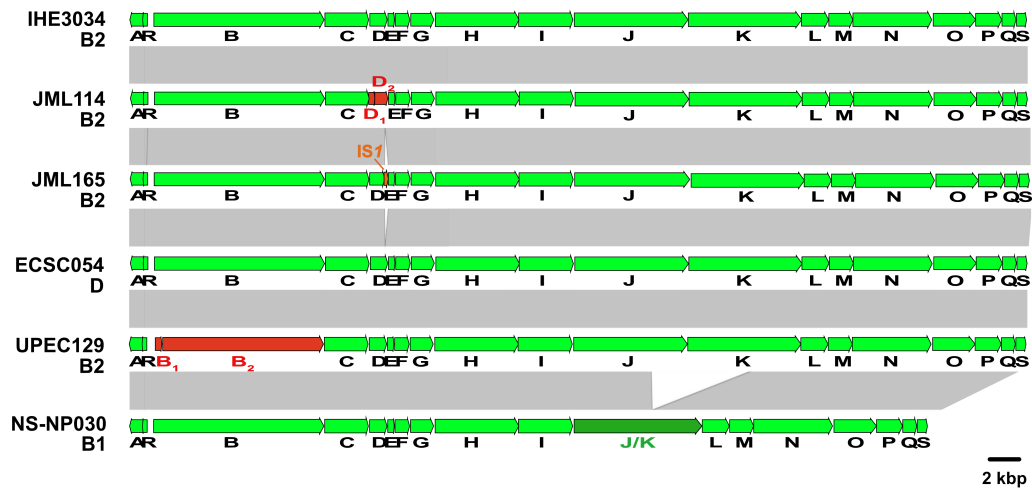
564  
 565  
 566 **Figure 6.** Cophylogeny of *pks* sequences and *E. coli* host strains. A comparison generated with  
 567 Phytools of *E. coli* core gene-based ML tree and *pks*-based NJ tree is shown, including the links  
 568 between *pks* and host strains (dashed lines). The phylogenetic groups (A, B1, B2 and D) are  
 569 indicated. The strains belonging to the major clonal complexes (CCs) are shown with coloured  
 570 names, including those of CC12 (containing 4 ST12 strains, one ST961 strain [ECSC078] and  
 571 one ST5389 strain [ECSC078]), CC14 (containing one ST14 strain [ECSC010] and 12 ST550  
 572 strains), CC73 (containing only ST73 strains) and CC95 (containing only ST95 strains). Strains

573 from CC95 that belong to subgroups A and B (as defined by *fimH* typing) or to serotype O1:H1  
574 are boxed. 7 reference *pks*-positive *E. coli* strains (536, ABU83972, CFT073, Nissle 1917,  
575 UTI89, IHE3034 and SP15) were included in both trees, whereas the non-*E. coli* strains  
576 carrying a *pks* island (i.e. *C. koseri* ATCC BAA-895, *E. aerogenes* EA1509E, *K. pneumoniae*  
577 1084, and *S. marcescens* AS012490) were included only in the *pks*-based tree.

578 ***The functionality of the cluster of genes of the pks island is conserved in the majority of the***  
579 ***enterobacterial strains***

580 We next investigated the functionality of the *pks* islands in *E. coli* strains belonging to various  
581 phylogroups and carrying phylogenetically distinct *pks* sequences, as well as in the *E. oleae*  
582 strain DAPP-PG531. The production of the genotoxin colibactin was directly investigated  
583 through the formation of DNA interstrand cross-links (ICLs) (Fig. S3A). The vast majority of  
584 the *E. coli* strains carrying the *pks* island (i.e. 83.5%) produced ICLs (Fig. 1). DNA-crosslinking  
585 was also observed for the *E. oleae* strain, and it was abrogated by adding purified colibactin  
586 self-resistance protein ClbS (Fig. S4), confirming the production of a *bona fide* colibactin by  
587 this strain carrying a less conserved sequence of the *pks* island. The *E. coli* genotoxic strains  
588 belonged to phylogroups B2, B1 and A (Fig. 1). Eighteen (16.5%) *pks*-positive *E. coli* isolates  
589 lacked a detectable interstrand crosslinking activity, including 15 strains from phylogroup B2,  
590 2 strains from phylogroup B1 and the single *pks*-positive strain from phylogroup D (Table 3).  
591 These strains did not cluster together in the core genome phylogenetic tree but instead were  
592 intertwined among genotoxic strains (Fig. 1, Fig. 5). To confirm the absence of genotoxicity,  
593 we tested the ability of these ICL-negative strains to trigger megalocytosis in cultured HeLa  
594 cells (Fig. S3B) and phosphorylation of histone H2AX (Fig. S3C), a robust marker for DNA  
595 damage in eukaryotic cells. To avoid cell lysis during infection, we assessed only non-  
596 hemolytic strains. No megalocytosis and no p-H2AX foci were detected in HeLa cells exposed  
597 to subsets of ICL-negative strains (Table 3; n=14 and n=5, respectively), even at a high  
598 multiplicity of infection, confirming the deficiency of these strains in colibactin production.  
599 These results showed that except for a few strains, *E. coli* strains carrying a *pks* island are

600 overwhelmingly capable of producing the genotoxin colibactin, regardless of their phylogenies  
601 and genomic configurations of their *pks* islands.  
602 To examine the reasons for the lack of genotoxic activity of the 18 ICL-negative *E. coli* strains,  
603 we further analyzed the sequence of the *pks* island from those strains. We identified genetic  
604 alterations of the *pks* island in 16 out of the 18 non-genotoxic isolates. Strain JML114 carried  
605 a single nucleotide deletion in *clbD* at position 172 (A), leading to the segregation of *clbD* into  
606 two ORFs (Fig. 7). Strain JML165 carried an *IS1* inserted at the 3'-end of *clbD* after position  
607 838. In strains UPEC129 and JML201, *clbB* was segregated into two ORFs due to nucleotide  
608 substitutions at positions 452 and 453 (AC to GA) (UPEC129; Fig. 7), and at position 872 (G  
609 to A) (JML201; data not shown), respectively. The genetic alterations identified in these four  
610 non-genotoxic strains each resulted in premature stop codons in *clbB* or *clbD* genes coding  
611 enzymes that are essential for the production of colibactin.



612

613 **Figure 7.** Comparison of the *pks* island sequence from the *E. coli* reference strain IHE3034  
614 with that of a selection of *pks*-positive but non-genotoxic *E. coli* isolates. Nucleotide sequence  
615 similarity (>99%) between different DNA regions is indicated by grey areas between the  
616 corresponding genetic maps. Fusion of two adjacent ORFs resulting from the deletion of a  
617 sequence overlapping the two ORFs is indicated in dark green. Adjacent ORF sequences  
618 resulting from the segregation of an original ORF following an insertion or deletion event are  
619 indicated in red. The *IS1* located in the *pks* island of strain JML165 is represented in orange.

620 Twelve ICL-negative strains carried a 5,651-bp deletion resulting in a *clbJK* fusion gene, as  
621 shown for strain NS-NP030 in Fig. 7. This deletion presumably resulted from recombination  
622 between two copies of a 1,480-bp homologous sequence located in *clbJ* and *clbK*. A PCR  
623 analysis of the corresponding region in the 109 *pks*-positive strains confirmed the presence of  
624 this deletion in the 12 strains, whereas the other 97 *pks*-positive strains contained full-length  
625 *clbJ* and *clbK* genes (Fig. 5). The 12 strains carrying the *clbJK* fusion were detected  
626 sporadically in the core genome phylogenetic tree (Fig. 1) and *pks* phylogenetic tree (Fig. 5),  
627 suggesting that occurrence of the deletion between *clbJ* and *clbK* arose from accidental  
628 recombination events. The predicted 2,440-amino acid hybrid ClbJK protein encoded by the  
629 *clbJK* fusion gene lacks the PKS module of ClbK necessary for the formation of stable cross-  
630 links (19). In agreement with this, the strains carrying this fusion were devoid of interstrand  
631 crosslinking activity and did not trigger megalocytosis nor histone H2AX phosphorylation in  
632 infected eukaryotic cells (Table 3). It was reported however that rat *E. coli* isolates carrying a  
633 *clbJK* fusion gene caused DNA damage or displayed cytotoxicity to HeLa cells (39, 40). This  
634 discrepancy with our results could be due to the use of distinct experimental conditions. The  
635 possibility that these rat isolates might produce additional genotoxins that would mask any  
636 colibactin deficiency caused by the *clbJK* fusion cannot be excluded either. Caution should also  
637 be observed during the assembly of sequencing reads as errors including deletions may be  
638 caused by the presence of tandem repeats.

639 For two other non-genotoxic isolates (CM1 and ECSC054), no mutation disrupting the *pks*  
640 genes were identified (Fig. 7; data not shown), suggesting that mutations located outside the  
641 *pks* island could negatively impact its expression. To test this hypothesis, we used plasmids  
642 pASK-clbR and pBAD-clbR, both overexpressing the *pks* regulator ClbR, and introduced either  
643 of them into the strain CM1 which was susceptible to antibiotic, in contrast to ECSC054. In the  
644 resulting CM1 transformants, colibactin activity was restored as seen by the formation of DNA

645 ICLs (data not shown). Thus, in this strain, the lack of genotoxic activity likely resulted from a  
646 negative regulation of the *pks* island through an unknown mechanism.

647 Functionality of the *pks* island was also examined through analysis of the lipid metabolite  
648 profiles of selected genotoxic (n=3) and non-genotoxic (n=8) *pks*-positive strains, and in  
649 particular for production of C14-Asn which was used as an indicator of the activity of the *pks*  
650 biosynthesis machinery. This lipopeptide is synthesized during the initial step of the  
651 biosynthesis process involving ClbN and ClbB, prior to elongation and final cleavage through  
652 the involvement of ClbC-H-I-J-K and ClbP, respectively (17). The production of C14-Asn was  
653 detected in all of the ICL-positive strains examined, SI-NP020, JML285 and KS-NP019 (*ca.*  
654 650-1200 pg/10e8 CFU) but not in the ICL-negative strain UPEC129 mutated in *clbB* (Fig. S5,  
655 Table 3). Interestingly, C14-Asn was detected (*ca.* 400-600 pg/10e8 CFU) in three ICL-  
656 negative strains carrying a *clbJK* fusion gene (SI-NP032, JML296 and NS-NP030) and in two  
657 ICL-negative strains carrying a mutated *clbD* gene (JML114 and JML165). The two non-  
658 genotoxic strains carrying intact *clb* genes (ECSC054 and CM1) produced either a very low  
659 level or no detectable C14-Asn, respectively (Fig. S5, Table 3). For the strain CM1 transformed  
660 with either plasmid pASK-clbR or pBAD-clbR (see above), overexpression of ClbR restored  
661 the production of C14-Asn (Fig. S5). These results suggest that even when the *pks* island does  
662 not allow production of active colibactin, enzymes from the *pks* pathway still produce  
663 metabolites with potential biological activities.

664

## 665 **DISCUSSION**

666 The acquisition of the *pks* island in the population of *E. coli* appears to have involved two  
667 distinct mechanisms differing by the presence or absence of a phage-type integrase. The  
668 integrase-mediated *pks* insertion pathway occurred mainly in B2 strains and resulted in *pks*

669 insertion into either of three *asn* tRNA genes (i.e. *asnU*, *asnV* or *asnW*). This potential for  
670 integration into several DNA targets is consistent with the observed conservation and genetic  
671 integrity of the *pks* integrative module, i.e. the integrase gene and the two direct repeats flanking  
672 the island. The flexibility of *pks* insertion is reminiscent of what has been described for the HPI  
673 of *Yersinia pseudotuberculosis* which is also able to insert into either of the three *Y.*  
674 *pseudotuberculosis* *asn* tRNA genes (41), in contrast to the immobile truncated form of the HPI  
675 in *E. coli* whose right direct repeat is deleted and whose location is fixed at the *asnT* tRNA gene  
676 (42). A divergent integrase sequence was found for the *pks* island inserted into the *asnU* tRNA  
677 gene compared to those inserted into the *asnV* or *asnW* tRNA gene. As the three *asn* tRNA  
678 sequences are 100% identical, the use of either of them as attachment site by slightly different  
679 integrases likely reflects distinct histories of *pks* acquisition. After *pks* chromosomal  
680 integration, the endogenous *pks* integrase promoter is replaced by the promoter of the upstream  
681 *asn* tRNA gene (Fig. S6), a configuration similar to that found for the HPI integrase promoter  
682 (43). Whether the site of integration influences the expression of the *pks* integrase and hence  
683 *pks* stability at the distinct *asn* tRNA *loci* is not known. In contrast to the integrase-mediated  
684 pathway, the *pks* chromosomal integration process in the B1 and A *E. coli* strains remains  
685 unclear as no site-specific recombinase-encoding gene was found near *pks* and chromosomal  
686 insertion occurred into a non-tRNA *locus*. In these strains, *pks* integration could have involved  
687 the participation of IS elements such as the IS66 whose truncated or intact copies were found  
688 to flank the *pks* island.

689 The cophylogeny analysis between the core genome- and *pks*-based phylogenetic trees shed  
690 further lights on the *pks* acquisition scenarios. In the case of the typical *pks*-positive B2 strains  
691 belonging to lineages from major CCs (i.e. CC12, CC14, CC73 and CC95), the congruence  
692 observed between both trees suggested that the *pks* island was horizontally acquired by the  
693 MRCA of these lineages, or by the MRCA of each of these, and then stably maintained in their

694 descendants through vertical transmission. The fact that strains of certain CC95 subgroups lack  
695 the *pks* island likely suggests that *pks* was lost during the evolution of these sublineages. Such  
696 a loss might be closely linked to the change in the relative fitness of CC95 subgroups underlying  
697 the variations observed in their spatial and temporal distribution in several continents (35). In  
698 the case of B1 or A strains, *pks* acquisition and dissemination likely occurred through sporadic  
699 lateral transfer events, as *pks*-positive B1 or A strains were scarce and not genetically related.  
700 The horizontal transferability of the *pks* island has previously been demonstrated using an *in*  
701 *vitro* approach where *pks* could be transferred together with the HPI via F' plasmid-mediated  
702 conjugation from a donor to a recipient *E. coli* strain (44). We propose that *pks* acquisition by  
703 the single *pks*-positive D strain ECSC054 was mediated by HGT, presumably from a B2 donor  
704 strain given the *pks* sequence relatedness observed between the D and B2 strains. HGT was  
705 also likely involved in the exchange of *pks* island between the three atypical B2 ST73 or ST95  
706 strains and a (yet unknown) phylogenetically distant donor strain, since their *pks* sequences did  
707 not cluster with those from other B2 ST73 or ST95 strains. This hypothesis was further  
708 supported by the identification, in these three isolates, of an ICE-like element inserted in their  
709 *pks* island. Similar ICE-like elements have previously been identified in three B1 *E. coli* isolates  
710 and other members of the *Enterobacteriaceae* such as *C. koseri*, *E. aerogenes* and *K.*  
711 *pneumoniae* (6). They could therefore play a role in *pks* dissemination in enterobacteria, as  
712 proposed for the self-transmissible ICE linked to the HPI identified in the *E. coli* strain ECOR31  
713 (38). Due to its lack of a complete DNA mobilization region (region II), we assume that the  
714 ICE-linked *pks* island is not self-transferable anymore. It might nevertheless correspond to a  
715 remnant of an ancient, complete and self-transmissible ICE-linked *pks* island that could have  
716 behaved as a large complex ICE and spread in enterobacteria before undergoing partial or entire  
717 deletion of the ICE region. To date, no bacterial strain carrying a complete ICE linked to *pks*  
718 has been identified and the origin of the *pks* island remains therefore elusive.



719 The ecological niche and/or genetic background of the bacterial strains probably had an impact  
720 on the acquisition and stable maintenance of *pks*. The high concentration of *pks*-positive strains  
721 in some CCs of the B2 group such as CC73 and CC95 suggests that *pks* might have contributed  
722 to their ecological and evolutionary success. CC73 and CC95 exhibit a similar phylogenetic  
723 history and are major ExPEC lineages, especially prior to the year 2000 where they were the  
724 most commonly detected (1, 45). They are persistent intestinal colonizers and successful extra-  
725 intestinal pathogens with the particularity of exhibiting lower multidrug resistance levels  
726 compared to other ExPEC lineages. By contrast, as our collection contained *E. coli* B2 strains  
727 from STs other than ST73 and ST95, it was interesting to note that *pks* was absent from STs  
728 corresponding to separate B2 lineages in the *E. coli* phylogenetic tree, including the ST131  
729 clonal complex which is associated with multidrug resistance and is now the most  
730 predominantly isolated ExPEC lineage worldwide (45). Consistent with the hypothesis of a *pks*  
731 acquisition by the MRCA(s) of CC73 and CC95 mentioned above, this finding suggests that  
732 such acquisition likely occurred after they diverged from the MRCA of CC131, i.e. before going  
733 through distinct evolutionary trajectories. The inversion of the upstream *asnW-asnU-asnV*  
734 tRNA-containing region which likely accompanied *pks* insertion into the *asnW* tRNA gene in  
735 the B2 group might have contributed to *pks* stabilization at this locus. Since the various *pks*-  
736 positive and -negative B2 lineages occupy the same ecological niche (i.e. primarily the  
737 intestinal tract of humans and animals), horizontal transfer of *pks* between them could have  
738 been expected, at least to some extent, which however was not revealed here. Several  
739 hypotheses can be proposed to explain this. First, some barriers to HGT might exist between  
740 members of distinct CCs, such as restriction-modification systems (46). Second, *pks* might have  
741 been transferred to recipient strains without providing adaptive value, thus resulting in its rapid  
742 loss. Third, as a crosstalk between virulence determinants and the chromosome backbone is  
743 required for the emergence of virulent clones (1), a specific chromosomal phylogenetic

744 background might be required for appropriate *pks* expression and production of an adaptive  
745 value, thereby constituting a prerequisite to the stable maintenance of the island.

746 The structure of the *pks* island is very well conserved among the *E. coli* population, with more  
747 than 99 % identity, suggesting that its integrity remains under strong structural and functional  
748 evolutionary constraints. We can speculate that, transcription and translation of the 19 *pks* genes  
749 of this 54 kb-long genomic island would be too high for the bacterial strains if the *pks* island  
750 did not bring a selective advantage to them. This is reinforced by the fact that only 18 out of  
751 109 *pks*-positive strains lacked genotoxic activity. The importance of *pks* biological role is  
752 highlighted by the numerous activities associated with this genetic island, including  
753 genotoxicity, anti-inflammatory activity, antibiotic and analgesic effects. Given its interplay  
754 with siderophores (enterobactin, salmochelin and yersiniabactin) and siderophore-microcins  
755 (MccM and MccH47) (10, 22, 47), the *pks* island contributes to bacterial competition through  
756 the acquisition of iron or the production of inhibitory compounds, respectively. Protection of  
757 bacterial cells from genomic degradation through the production of ClbS could also be  
758 advantageous to *pks*-carrying strains, as this multifunctional protein not only directly  
759 inactivates colibactin but also protects bacterial DNA from nucleolytic degradation by  
760 nucleases (48). We also observed that non-genotoxic *E. coli* strains carrying an altered *pks*  
761 island still produced the prodrug motif C14-Asn synthesized at the early stage of the  
762 biosynthesis process, suggesting that yet-to-be-discovered bioactive compounds are produced  
763 by these strains. Given the high conservation observed for the *pks* island in *E. coli*, we can thus  
764 speculate that colibactin is a very important genotoxin but that *pks*-derived synthesis of other  
765 secondary metabolites could also be an advantage for *E. coli*.

766 Although our collection is characterized by a large diversity of *E. coli* strains from various  
767 phylogenetic groups and STs, one limitation of this study is that only strains from Japan were

768 included, which may not be representative of *pks* distribution in a global collection of *E. coli*  
769 isolates from worldwide sources.

770 In conclusion, the various genetic configurations of the *pks* island and its distribution in the *E.*  
771 *coli* phylogenetic tree imply the existence of various scenarios for the introduction and spread  
772 of *pks* into the *E. coli* population. The presence of a functional *pks* island was demonstrated for  
773 the majority of the *pks*-positive strains, suggesting that the *pks* island is under selective pressure  
774 for the adaptation of *E. coli* to various ecological niches, through the production of colibactin  
775 or other secondary metabolites.

776

## 777 **AUTHOR STATEMENTS**

### 778 **Author contributions**

779 Conceptualisation: FA, TH, PB, YO, EO. Methodology: FA, TH, PB, YO, EO. Validation: FA,  
780 TH, PB, YO, EO. Formal Analysis: FA, AP, CC, YA, TH, PB, YO, EO. Investigation: FA, AP,  
781 YA, CC, NBG, CM, JPN, PB, YO, EO. Resources: TH, YO, EO. Data Curation: FA, AP, YA,  
782 TH, PB, YO, EO. Writing – Original Draft: FA, TH, YO, EO. Writing - Review and Editing:  
783 FA, AP, CC, HB, JPN, TH, PB, YO, EO. Visualization: FA, YA, CC, JPN, PB, YO, EO.  
784 Supervision: FA, HB, TH, YO, EO. Project Administration: FA, EO. Funding Acquisition: TH,  
785 YO, EO.

### 786 **Conflicts of interest**

787 The authors declare that there are no conflicts of interest.

### 788 **Funding information**

789 This work was supported by fundings from the National French Institute of Health and Medical  
790 Research (INSERM) to Camille Chagneau and the région Occitanie (grant ALDOCT-000610)  
791 and Ministère de l'Agriculture to Alexandre Perrat.

792 **Acknowledgments**

793 We thank Claire Hoede and Sarah Maman (SIGENAE group) and the GENOTOUL  
794 bioinformatics platform for providing computational resources. We also thank Pauline Le  
795 Faouder and the METATOUL lipidomic platform for their support in the analysis of the lipid  
796 metabolite profiles.

797 **REFERENCES**

- 798 1. Denamur E, Clermont O, Bonacorsi S, Gordon D. The population genetics of pathogenic  
799 *Escherichia coli*. *Nature reviews Microbiology*. 2020.
- 800 2. Tenailon O, Skurnik D, Picard B, Denamur E. The population genetics of commensal  
801 *Escherichia coli*. *Nature reviews Microbiology*. 2010;8(3):207-17.
- 802 3. Ochman H, Lawrence JG, Groisman EA. Lateral gene transfer and the nature of  
803 bacterial innovation. *Nature*. 2000;405(6784):299-304.
- 804 4. Kaper JB, Nataro JP, Mobley HL. Pathogenic *Escherichia coli*. *Nature reviews*  
805 *Microbiology*. 2004;2(2):123-40.
- 806 5. Nougayrède JP, Homburg S, Taieb F, Boury M, Brzuszkiewicz E, Gottschalk G, et al.  
807 *Escherichia coli* induces DNA double-strand breaks in eukaryotic cells. *Science (New York, NY)*.  
808 2006;313(5788):848-51.
- 809 6. Putze J, Hennequin C, Nougayrède JP, Zhang W, Homburg S, Karch H, et al. Genetic  
810 structure and distribution of the colibactin genomic island among members of the family  
811 *Enterobacteriaceae*. *Infection and immunity*. 2009;77(11):4696-703.
- 812 7. Engel P, Vizcaino MI, Crawford JM. Gut symbionts from distinct hosts exhibit genotoxic  
813 activity via divergent colibactin biosynthesis pathways. *Applied and environmental*  
814 *microbiology*. 2015;81(4):1502-12.
- 815 8. Bondarev V, Richter M, Romano S, Piel J, Schwedt A, Schulz-Vogt HN. The genus  
816 *Pseudovibrio* contains metabolically versatile bacteria adapted for symbiosis. *Environmental*  
817 *microbiology*. 2013;15(7):2095-113.
- 818 9. Marcq I, Martin P, Payros D, Cuevas-Ramos G, Boury M, Watrin C, et al. The genotoxin  
819 colibactin exacerbates lymphopenia and decreases survival rate in mice infected with  
820 septicemic *Escherichia coli*. *The Journal of infectious diseases*. 2014;210(2):285-94.
- 821 10. Martin P, Marcq I, Magistro G, Penary M, Garcie C, Payros D, et al. Interplay between  
822 siderophores and colibactin genotoxin biosynthetic pathways in *Escherichia coli*. *PLoS*  
823 *pathogens*. 2013;9(7):e1003437.
- 824 11. McCarthy AJ, Martin P, Cloup E, Stabler RA, Oswald E, Taylor PW. The Genotoxin  
825 Colibactin Is a Determinant of Virulence in *Escherichia coli* K1 Experimental Neonatal Systemic  
826 Infection. *Infection and immunity*. 2015;83(9):3704-11.
- 827 12. Arthur JC, Perez-Chanona E, Mühlbauer M, Tomkovich S, Uronis JM, Fan TJ, et al.  
828 Intestinal inflammation targets cancer-inducing activity of the microbiota. *Science (New York,*  
829 *NY)*. 2012;338(6103):120-3.
- 830 13. Cougnoux A, Dalmasso G, Martinez R, Buc E, Delmas J, Gibold L, et al. Bacterial  
831 genotoxin colibactin promotes colon tumour growth by inducing a senescence-associated  
832 secretory phenotype. *Gut*. 2014;63(12):1932-42.

- 833 14. Cuevas-Ramos G, Petit CR, Marcq I, Boury M, Oswald E, Nougayrède JP. *Escherichia coli*  
834 induces DNA damage in vivo and triggers genomic instability in mammalian cells. *Proceedings*  
835 *of the National Academy of Sciences of the United States of America*. 2010;107(25):11537-42.
- 836 15. Bossuet-Greif N, Vignard J, Taieb F, Mirey G, Dubois D, Petit C, et al. The Colibactin  
837 Genotoxin Generates DNA Interstrand Cross-Links in Infected Cells. *mBio*. 2018;9(2).
- 838 16. Taieb F, Petit C, Nougayrède JP, Oswald E. The Enterobacterial Genotoxins: Cytolethal  
839 Distending Toxin and Colibactin. *EcoSal Plus*. 2016;7(1).
- 840 17. Brotherton CA, Balskus EP. A prodrug resistance mechanism is involved in colibactin  
841 biosynthesis and cytotoxicity. *Journal of the American Chemical Society*. 2013;135(9):3359-  
842 62.
- 843 18. Wallenstein A, Rehm N, Brinkmann M, Selle M, Bossuet-Greif N, Sauer D, et al. ClbR Is  
844 the Key Transcriptional Activator of Colibactin Gene Expression in *Escherichia coli*. *mSphere*.  
845 2020;5(4).
- 846 19. Shine EE, Xue M, Patel JR, Healy AR, Surovtseva YV, Herzon SB, et al. Model Colibactins  
847 Exhibit Human Cell Genotoxicity in the Absence of Host Bacteria. *ACS chemical biology*.  
848 2018;13(12):3286-93.
- 849 20. Vizcaino MI, Engel P, Trautman E, Crawford JM. Comparative metabolomics and  
850 structural characterizations illuminate colibactin pathway-dependent small molecules.  
851 *Journal of the American Chemical Society*. 2014;136(26):9244-7.
- 852 21. Pérez-Berezo T, Pujo J, Martin P, Le Faouder P, Galano JM, Guy A, et al. Identification  
853 of an analgesic lipopeptide produced by the probiotic *Escherichia coli* strain Nissle 1917.  
854 *Nature communications*. 2017;8(1):1314.
- 855 22. Massip C, Branchu P, Bossuet-Greif N, Chagneau CV, Gaillard D, Martin P, et al.  
856 Deciphering the interplay between the genotoxic and probiotic activities of *Escherichia coli*  
857 Nissle 1917. *PLoS pathogens*. 2019;15(9):e1008029.
- 858 23. Dubois D, Delmas J, Cady A, Robin F, Sivignon A, Oswald E, et al. Cyclomodulins in  
859 urosepsis strains of *Escherichia coli*. *Journal of clinical microbiology*. 2010;48(6):2122-9.
- 860 24. Johnson JR, Johnston B, Kuskowski MA, Nougayrède JP, Oswald E. Molecular  
861 epidemiology and phylogenetic distribution of the *Escherichia coli* pks genomic island. *Journal*  
862 *of clinical microbiology*. 2008;46(12):3906-11.
- 863 25. Arimizu Y, Kirino Y, Sato MP, Uno K, Sato T, Gotoh Y, et al. Large-scale genome analysis  
864 of bovine commensal *Escherichia coli* reveals that bovine-adapted *E. coli* lineages are serving  
865 as evolutionary sources of the emergence of human intestinal pathogenic strains. *Genome*  
866 *research*. 2019;29(9):1495-505.
- 867 26. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies  
868 from short and long sequencing reads. *PLoS computational biology*. 2017;13(6):e1005595.
- 869 27. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, et al. Roary: rapid large-  
870 scale prokaryote pan genome analysis. *Bioinformatics (Oxford, England)*. 2015;31(22):3691-3.
- 871 28. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, et al. SNP-sites: rapid  
872 efficient extraction of SNPs from multi-FASTA alignments. *Microbial genomics*.  
873 2016;2(4):e000056.
- 874 29. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with  
875 thousands of taxa and mixed models. *Bioinformatics (Oxford, England)*. 2006;22(21):2688-90.
- 876 30. Letunic I, Bork P. Interactive tree of life (iTOL) v3: an online tool for the display and  
877 annotation of phylogenetic and other trees. *Nucleic acids research*. 2016;44(W1):W242-5.
- 878 31. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile  
879 and open software for comparing large genomes. *Genome biology*. 2004;5(2):R12.

- 880 32. Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis  
881 Version 7.0 for Bigger Datasets. *Molecular biology and evolution*. 2016;33(7):1870-4.
- 882 33. Revell LJ. Phytools: an R package for phylogenetic comparative biology (and other  
883 things). *Methods in Ecology and Evolution*. 2012;3:217-23.
- 884 34. Carver TJ, Rutherford KM, Berriman M, Rajandream MA, Barrell BG, Parkhill J. ACT: the  
885 Artemis Comparison Tool. *Bioinformatics (Oxford, England)*. 2005;21(16):3422-3.
- 886 35. Gordon DM, Geyik S, Clermont O, O'Brien CL, Huang S, Abayasekara C, et al. Fine-Scale  
887 Structure Analysis Shows Epidemic Patterns of Clonal Complex 95, a Cosmopolitan *Escherichia*  
888 *coli* Lineage Responsible for Extraintestinal Infection. *mSphere*. 2017;2(3).
- 889 36. Bossuet-Greif N, Belloy M, Boury M, Oswald E, Nougayrede J-P. Protocol for HeLa Cells  
890 Infection with *Escherichia coli* Strains Producing Colibactin and Quantification of the Induced  
891 DNA-damage. *Bio-protocol*. 2017;7(16):e2520.
- 892 37. Tronnet S, Oswald E. Quantification of Colibactin-associated Genotoxicity in HeLa Cells  
893 by In Cell Western (ICW) Using  $\gamma$ -H2AX as a Marker. *Bio-protocol*. 2018;8(6):e2771.
- 894 38. Schubert S, Dufke S, Sorsa J, Heesemann J. A novel integrative and conjugative element  
895 (ICE) of *Escherichia coli*: the putative progenitor of the *Yersinia* high-pathogenicity island.  
896 *Molecular microbiology*. 2004;51(3):837-48.
- 897 39. Fabian NJ, Mannion AJ, Feng Y, Madden CM, Fox JG. Intestinal colonization of  
898 genotoxic *Escherichia coli* strains encoding colibactin and cytotoxic necrotizing factor in small  
899 mammal pets. *Veterinary microbiology*. 2020;240:108506.
- 900 40. Kurnick SA, Mannion AJ, Feng Y, Madden CM, Chamberlain P, Fox JG. Genotoxic  
901 *Escherichia coli* Strains Encoding Colibactin, Cytolethal Distending Toxin, and Cytotoxic  
902 Necrotizing Factor in Laboratory Rats. *Comparative medicine*. 2019;69(2):103-13.
- 903 41. Buchrieser C, Brosch R, Bach S, Guiyoule A, Carniel E. The high-pathogenicity island of  
904 *Yersinia pseudotuberculosis* can be inserted into any of the three chromosomal *asn* tRNA  
905 genes. *Molecular microbiology*. 1998;30(5):965-78.
- 906 42. Schubert S, Darlu P, Clermont O, Wieser A, Magistro G, Hoffmann C, et al. Role of  
907 intraspecies recombination in the spread of pathogenicity islands within the *Escherichia coli*  
908 species. *PLoS pathogens*. 2009;5(1):e1000257.
- 909 43. Rakin A, Noelting C, Schropp P, Heesemann J. Integrative module of the high-  
910 pathogenicity island of *Yersinia*. *Molecular microbiology*. 2001;39(2):407-15.
- 911 44. Messerer M, Fischer W, Schubert S. Investigation of horizontal gene transfer of  
912 pathogenicity islands in *Escherichia coli* using next-generation sequencing. *PloS one*.  
913 2017;12(7):e0179880.
- 914 45. Manges AR, Geum HM, Guo A, Edens TJ, Fibke CD, Pitout JDD. Global Extraintestinal  
915 Pathogenic *Escherichia coli* (ExPEC) Lineages. *Clinical microbiology reviews*. 2019;32(3).
- 916 46. Oliveira PH, Touchon M, Rocha EP. Regulation of genetic flux between bacteria by  
917 restriction-modification systems. *Proceedings of the National Academy of Sciences of the*  
918 *United States of America*. 2016;113(20):5658-63.
- 919 47. Massip C, Chagneau CV, Boury M, Oswald E. The synergistic triad between microcin,  
920 colibactin, and salmochelin gene clusters in uropathogenic *Escherichia coli*. *Microbes and*  
921 *infection*. 2020;22(3):144-7.
- 922 48. Molan K, Podlesek Z, Hodnik V, Butala M, Oswald E, Žgur Bertok D. The *Escherichia coli*  
923 colibactin resistance protein *ClbS* is a novel DNA binding protein that protects DNA from  
924 nucleolytic degradation. *DNA repair*. 2019;79:50-4.
- 925