

# Diversity and prevalence of colibactin- and yersiniabactin encoding mobile genetic elements in enterobacterial populations: insights into evolution and co-existence of two bacterial secondary metabolite determinants

Haleluya Wami<sup>1#</sup>, Alexander Wallenstein<sup>1#</sup>, Daniel Sauer<sup>2</sup>, Monika Stoll<sup>3</sup>, Rudolf von Büнау<sup>4</sup>, Eric Oswald<sup>5</sup>, Rolf Müller<sup>2</sup>, Ulrich Dobrindt<sup>1\*</sup>

<sup>1</sup>Institute of Hygiene, University of Münster, Münster, Germany

<sup>2</sup>Department of Microbial Natural Products, Helmholtz Institute for Pharmaceutical Research Saarland, Helmholtz Center for Infection Research, Saarland University, Campus E8 1, Saarbrücken, Germany

<sup>3</sup>Department of Genetic Epidemiology, Institute of Human Genetics, University of Münster, Münster, Germany.

<sup>4</sup>Ardeypharm GmbH, Herdecke, Germany

<sup>5</sup>IRSD, Université de Toulouse, INSERM, INRA, ENVT, UPS, Toulouse, France

\*For correspondence:

Prof. Dr. Ulrich Dobrindt  
Institute of Hygiene  
Mendelstraße 7  
48149 Münster  
Germany  
Tel. +49 (0)251 980 28775  
Fax: +49 (0)251 980 2868  
Email: [dobrindt@uni-muenster.de](mailto:dobrindt@uni-muenster.de)

Key words: High Pathogenicity Island, secondary metabolite, polyketide, cytopathic effect, *E. coli*, *Klebsiella*, *Citrobacter*

## 1 **1 Abstract**

2 The bacterial genotoxin colibactin interferes with the eukaryotic cell cycle by causing double-  
3 stranded DNA breaks. It has been linked to bacterially induced colorectal cancer in humans.  
4 Colibactin is encoded by a 54-kb genomic region in *Enterobacteriaceae*. The colibactin genes  
5 commonly co-occur with the yersiniabactin biosynthetic determinant. Investigating the  
6 prevalence and sequence diversity of the colibactin determinant and its linkage to the  
7 yersiniabactin operon in prokaryotic genomes, we discovered mainly species-specific lineages  
8 of the colibactin determinant and classified three main structural settings of the colibactin-  
9 yersiniabactin genomic region in *Enterobacteriaceae*. The colibactin gene cluster has a similar  
10 but not identical evolutionary track to that of the yersiniabactin operon. Both determinants  
11 could have been acquired on several occasions and/or exchanged independently between  
12 enterobacteria by horizontal gene transfer. Integrative and conjugative elements play(ed) a  
13 central role in the evolution and structural diversity of the colibactin-yersiniabactin genomic  
14 region. Addition of an activating and regulating module (*clbAR*) to the biosynthesis and  
15 transport module (*clbB-S*) represents the most recent step in the evolution of the colibactin  
16 determinant. In a first attempt to correlate colibactin expression with individual lineages of  
17 colibactin determinants and different bacterial genetic backgrounds, we compared colibactin  
18 expression of selected enterobacterial isolates *in vitro*. Colibactin production in the tested  
19 *Klebsiella* spp. and *Citrobacter koseri* strains was more homogeneous and generally higher  
20 than that in most of the *E. coli* isolates studied. Our results improve the understanding of the  
21 diversity of colibactin determinants and its expression level, and may contribute to risk  
22 assessment of colibactin-producing enterobacteria.

## 23 **2 Introduction**

24 The non-ribosomal peptide/polyketide hybrid colibactin is a secondary metabolite found in a  
25 variety of bacterial species of the family *Enterobacteriaceae*. The colibactin biosynthetic  
26 machinery is encoded by a 54-kb large polyketide synthase (*pks*) or *clb* genomic island[1],  
27 which includes 19 genes. The largest part of the island consists of a section of overlapping or  
28 closely spaced genes: *clbB* to *clbL* and *clbN* to *clbQ*, which are aligned on the same strand and  
29 code for components of the biosynthesis complex. The colibactin assembly line is  
30 supplemented with a dedicated transporter, encoded by *clbM*, and a resistance-conferring

31 protein encoded by *clbS* [2, 3]. Two additional genes required for colibactin production are  
32 located ca. 400 bp upstream of the first biosynthesis gene *clbB* in the opposing reading  
33 direction: the *clbR* gene coding for an auto-activating, *pks* island-specific transcription factor  
34 and the phosphopantetheinyl transferase-encoding gene *clbA*, which is crucial for activation  
35 of polyketide biosynthesis complexes (Fig. 1) [4-6]. Between these two divergent transcription  
36 units, there is a "variable number of tandem repeat" (VNTR) region, which comprises a varying  
37 number of a repeating octanucleotide sequence (5'-ACAGATAC-3') depending on the isolate  
38 [2].

39 Recently, the structure of the colibactin molecule has been proposed [3, 7-9]. Yet, the  
40 biological role of colibactin is still under discussion. Colibactin can interfere with the  
41 progression of the eukaryotic cell cycle, presumably by cross-linking DNA resulting in double-  
42 strand DNA breaks in genomic instability in eukaryotes [1, 10, 11]. The ability to produce  
43 colibactin has been described to increase the pathogenic potential of the producing bacteria  
44 and to promote colorectal cancer development [12-17], but has also been related to  
45 beneficial effects to the host [3, 18-20].

46 Initially, the *pks* island has been described in extraintestinal pathogenic *E. coli* to be  
47 chromosomally inserted into the *asnW* tRNA locus in close proximity to another tRNA(*Asn*)  
48 gene-associated pathogenicity island, the so-called "high pathogenicity island" (HPI). The HPI  
49 harbors an additional polyketide determinant coding for the metallophore yersiniabactin  
50 biosynthetic machinery [1, 21]. As members of the *Enterobacteriaceae* are generally not  
51 known as archetypal secondary metabolite producers the origin of the *pks* island remains to  
52 be further investigated. Interestingly, the colibactin determinant has also been detected as  
53 part of an "integrative and conjugative element" (ICE) in different enterobacteria. This ICE also  
54 integrates near a tRNA(*Asn*) locus into the bacterial chromosome and commonly carries the  
55 yersiniabactin gene cluster [2, 22]. It has been suggested that the close linkage observed  
56 between the colibactin and yersiniabactin gene clusters results from the functional  
57 interconnection between the colibactin and yersiniabactin biosynthetic pathways via the  
58 phosphopantetheinyl transferase ClbA, which can also contribute to the biosynthesis of  
59 yersiniabactin [5, 23]. The highly conserved colibactin determinant has so far been detected in  
60 a spectrum of strains belonging to the *Enterobacteriaceae* family: most commonly among *E.*  
61 *coli* strains of the phylogroup B2, followed by *Klebsiella pneumoniae* isolates, but also in  
62 *Citrobacter koseri* and *Klebsiella aerogenes* [2, 18, 24]. Less conserved variants or homologs of

63 the colibactin gene cluster have been described phenotypically or based on nucleotide  
64 sequence data in an *Erwinia oleae* strain, the honey bee symbiont *Frischella perrara*, and the  
65 marine  $\alpha$ -proteobacterium *Pseudovibrio* [24, 25].

66 Based on the low sequence similarity between the enterobacterial colibactin genes and the  
67 two homologous polyketide determinants in *F. perrara* and *Pseudovibrio* as well as on their  
68 association with mobile genetic elements (MGEs) or at least mobility-associated genes, one  
69 can hypothesize that the colibactin gene cluster is spread by horizontal gene transfer, maybe  
70 via an ICE-like element [2, 24]. While in most *Enterobacteriaceae* the colibactin determinant is  
71 typically associated with an ICE, the characteristic mobility and transfer features of an ICE are  
72 absent in the sequence context of the *pks* island in *E. coli* phylogroup B2 strains. Nevertheless,  
73 the *pks* island in *E. coli* remains mobilizable and transferable through external factors,  
74 supporting the hypothesis that former MGEs can undergo a stabilization (homing) process  
75 upon their chromosomal integration [26-28].

76 Studies addressing the prevalence of the colibactin genes were so far mainly focused on  
77 *Klebsiella* spp. or *Escherichia coli* backgrounds. The scarcity of data in other prokaryotic  
78 species regarding its distribution and the structure of the associated MGE makes it challenging  
79 to reliably further characterize the transmission and evolution of this polyketide determinant  
80 [18]. Previous data show that the prevalence varied from 5.3% to 25.6% in *Klebsiella* and from  
81 9.5% to 58% in *Escherichia*, highlighting an enrichment of *pks* island in specific ecological  
82 niches, whereas studies with a broader screening approach resulted in a prevalence of 14% in  
83 *Klebsiella* and 9.5% in *Escherichia* isolates, respectively [2, 29-34]. Notably, in health-related  
84 studies, a higher association of the *clb* genes was observed amongst strains with an increased  
85 virulence potential, with a prevalence of as much as 78.8% for *Klebsiella* subgroups and 72.7%  
86 for colorectal cancer-associated *E. coli* isolates [29, 32, 35-37]. The colibactin genes are  
87 frequently found in hyper-virulent and multidrug-resistant *K. pneumoniae* isolates [38, 39].

88 The obvious prevalence of the colibactin gene cluster in specific enterobacterial species  
89 combined with the description of more distantly related homologous determinants has  
90 sparked our interest in a better understanding of the spread and evolution of the colibactin  
91 determinant and its genetic context in bacteria. Therefore, our study aimed to investigate the  
92 prevalence and diversity of the colibactin determinant also in isolates outside of the  
93 *Enterobacteriaceae* family. Furthermore, we compared colibactin expression levels among  
94 enterobacterial isolates carrying different lineages of colibactin determinants as a first

95 attempt to assess the functional context of the bacterial genetic background, pathogenicity,  
96 and colibactin expression.

97

### 98 **3 Impact Statement**

99 Colibactin can act as a bacterial genotoxin and thus promote colorectal cancer development.  
100 Little is known about the origin, diversity and prevalence of the colibactin genes (*clb*) within  
101 prokaryotes. The *clb* genes are closely associated with pathogenicity islands or integrative and  
102 conjugative elements (ICEs). We screened roughly 375,000 prokaryotic genomes to analyze the  
103 diversity and evolution of such mobile genetic elements among bacterial populations.  
104 Interestingly, *clb* genes are only present in a subgroup of enterobacteria, mainly *E. coli*,  
105 *Klebsiella* spp. and *Citrobacter koseri*. The *clb* determinant, together with the yersiniabactin  
106 (*ybt*) gene cluster, belong to an ICE in most of the *clb*-positive enterobacteria, especially in  
107 *Klebsiella*. We show that both determinants, though in principle freely transferable within  
108 bacteria, have a mainly species-specific phylogeny, and that colibactin expression levels were  
109 species-independent. Recombination promoted the structural diversification of the ICE in  
110 different species, including its successive degeneration that led to the establishment of the  
111 colibactin and yersiniabactin islands in *E. coli* phylogroup B2 strains. Our results not only  
112 illustrate differing evolutionary tracks of the *clb* and *ybt* determinants in different  
113 enterobacterial species, but also highlight the important of ICEs for genomic variability in  
114 enterobacteria and the evolution of archetypal pathogenicity islands.

115

## 116 4 Methods

### 117 4.1 Bacterial strains and media

118 For cultivation, bacteria were grown as batch cultures in lysogeny broth (LB) (10g/l tryptone,  
119 5g/l yeast extract, 5g/l sodium chloride) at 37°C. Strains used in this study are listed in the  
120 Supplemental Material (Table S1).

### 121 4.2 DNA extraction and sequencing

122 DNA extraction of the enterobacterial strains was performed using the MagAttract® HMW  
123 DNA Kit (Qiagen, Hilden, Germany) according to the manufacturer's recommendations. To  
124 prepare paired-end libraries we used the Nextera XT DNA Library Preparation kit (Illumina,  
125 San Diego, CA, USA). Libraries were sequenced on the Illumina MiSeq sequencing platform  
126 using v2 sequencing chemistry (500 cycles) or on the Illumina NextSeq500 system using v2.5  
127 chemistry (300 cycles). Accession numbers of in-house sequences submitted to the NCBI  
128 GenBank database are included in Supplemental Material (Table S1).

### 129 4.3 Genome selection and phylogenetic analysis

130 All genome sequences not generated in this study were obtained from publically available  
131 prokaryotic genomes (NCBI GenBank). The quality of in-house sequenced genomes was  
132 checked with FastQC v0.11.5 (<https://github.com/chgibb/FastQC0.11.5/blob/master/fastqc>),  
133 and low-quality reads were trimmed using Sickle v1.33 (<https://github.com/najoshi/sickle>). The  
134 processed reads were *de novo* assembled with SPAdes v3.13.1 [40] and annotated with  
135 prokka v1.12 [41]. The genomes were screened for the presence of >45 kb of the complete  
136 *pks* genomic island using standalone BLAST+ v2.8.1 [42] and antiSMASH v5.0.0 [43]. The *pks*  
137 island found in the genome of the *E. coli* strain M1/5 (accession no. [CP053296](#)) was used as a  
138 reference sequence. The VNTR and the sequence stretch of the *pks* island that spans between  
139 *clbJ* and *clbK* were excluded from analysis as these regions are prone to misassembly.

140 The contigs that align to the colibactin genes were ordered using ABACAS v1.3.1 [44] and  
141 multiple sequence alignment was generated using Kalign v3.1.1 [45]. Recombinant regions  
142 were detected and removed using Gubbins v2.4.1 [46]. The recombination filtered  
143 polymorphisms were then used to generate a maximum likelihood phylogeny of the colibactin  
144 determinant using RAxML v8.2.11 [47] under the GTR-GAMMAX model from 9974  
145 polymorphic sites. The branch support of the maximum likelihood tree was estimated by

146 bootstrap analysis of 200 replicate trees. The homologous gene cluster found in *F. perrara*  
147 was used as an outgroup. The phylogeny of the corresponding *ybt* islands was generated with  
148 a similar approach. The generated trees were visualized using itoL (<https://itol.embl.de>).

#### 149 **4.4 Phylo-grouping of *pks*-positive strains**

150 The *E. coli* and *Klebsiella* strains that harbored the colibactin gene cluster were allocated to  
151 their corresponding sequence types using mlst v2.16.1 (<https://github.com/tseemann/mlst>),  
152 which detects sequence types using the PubMLST typing schemes. The *Escherichia* strains  
153 were further classified into their phylogenetic lineages using the standalone tool, EzClermont  
154 v0.4.5 (<https://github.com/nickp60/EzClermont>).

155 The analysis of the diversity of the colibactin and yersiniabactin gene clusters involved  
156 virulence gene multi-locus sequence typing (MLST) for both polyketide determinants as  
157 previously described [38]. Briefly, the allele sequences of sixteen genes of the colibactin gene  
158 cluster (*clbACDEFGHILMNOPQR*) as well as of eleven genes of the yersiniabactin determinant  
159 (*fyuA*, *ybtE*, *ybtT*, *ybtU*, *irp1*, *irp2*, *ybtA*, *ybtP*, *ybtQ*, *ybtX*, *ybtS*) were extracted from the  
160 individual genomes and analyzed for allelic variations. Each observed combination of alleles  
161 was assigned a unique colibactin sequence type (CbST, listed in Table S5) or yersiniabactin  
162 sequence type (YbST, listed in Table S6).

#### 163 **4.5 Variable number tandem repeat (VNTR) detection**

164 The VNTR copy number present within the colibactin determinant (upstream of *clbR*) was  
165 detected using the standalone version of tandem repeats finder v4.09 [48]. The VNTR copy  
166 number distribution was visualized using R v3.4.3 (<https://www.r-project.org/index.html>).

#### 167 **4.6 Detection of *E. coli* virulence markers for pathotyping**

168 For pathotyping, the *clb*-positive *E. coli* strains were *in silico* screened for the presence of  
169 different *E. coli* pathotype marker genes using BLAST+ v2.8.1 (Supplemental Material, Tab. S).  
170 These genes were used as markers for *E. coli* pathotypes: enteroaggregative *E. coli* (EAEC),  
171 enterohemorrhagic *E. coli* (EHEC), enteropathogenic *E. coli* (EPEC), enterotoxigenic *E. coli*  
172 (ETEC), diffusely adhering *E. coli* (DAEC), uropathogenic *E. coli* (UPEC), and newborn  
173 meningitis-causing *E. coli* (NMEC).

#### 174 **4.7 Quantification of colibactin expression through N-myristoyl-D-asparagine**

175 Following an approach described by Bian and colleagues [49], a collection of colibactin-  
176 producing strains of the main species harboring the colibactin determinant was characterized  
177 for their ability to produce colibactin under *in vitro* growth conditions. For this purpose, we  
178 quantified N-myristoyl-D-asparagine (N-Myr-D-asparagine) a byproduct during colibactin  
179 maturation. The amount of this intermediate extrapolates the resulting colibactin amount  
180 produced. After growing the bacteria for 24 h at 37°C in glass tubes in 5 ml LB supplemented  
181 with 200 µl of a water/XAD-16-resin slurry, bacterial cells were harvested by centrifugation.  
182 The pelleted bacteria-slurry mixes were sedimented, filtered, and dissolved three times in  
183 acetone with increasing volume (12ml, 100 ml, and finally 200ml). Afterward, the solvent was  
184 exchanged by rotary evaporation and replaced by 1.6 ml methanol. The sample was further  
185 concentrated by centrifugation (10 min, 15000 rpm at 4°C), followed by drying 1,5 ml of the  
186 solution in a vacuum centrifuge and subsequent resuspension in 50 µl methanol. 30 µl of  
187 these processed samples were measured by Ultra Performance Liquid Chromatography  
188 coupled to High-Resolution Mass Spectrometry (UPLC-HRMS) conducted on a Thermo  
189 Scientific Ultimate 3000 RS with a Waters Acquity BEH 100\*2.1mm 1.7µm 130A column  
190 (eluent A: 0.1% formic acid in ddH<sub>2</sub>O, eluent B: 0.1% formic acid in acetonitrile), where a flow  
191 rate of 0.6ml/min followed by a Bruker Maxis II-4G, 150-2500 m/z and a scan rate of 2Hz was  
192 applied. To enable quantification of N-Myr-D-asparagine, we used 250 mM cinnarizine as  
193 internal standard and normalized peak areas based on the internal standard and the optical  
194 density (OD<sub>600</sub>) of the bacterial culture.

195



## 196 5 Results

### 197 5.1 Prevalence of colibactin determinant

198 Of the 374,754 publicly accessible prokaryotic genomes (as of 30.06.2019) that were screened  
199 for the presence of the colibactin gene cluster, 1,969 genomes carried this polyketide-  
200 encoding operon. An additional 200 *clb*-positive enterobacterial genomes determined in-  
201 house were added to the analysis (Tab. 1, Supplemental Material, Tab. S1). The *clb* gene  
202 cluster was detected in several enterobacterial species, most frequently in *E. coli* and *K.*  
203 *pneumoniae* isolates, but also to a lesser extent in *K. aerogenes*, *C. koseri*, *E. cloacae*, *E.*  
204 *hormaechei*, *K. michiganensis*, *S. marcescens*, and *E. oleae*. The colibactin determinant was,  
205 however, not detectable in 112,546 *Salmonella enterica* and 41 *S. bongori* genomes  
206 (Supplemental Material, Tab. S4), but in one out of eight genomes of unspecified *Salmonella*  
207 isolates. We did not detect the *clb* genes in 2,634 *Shigella* spp., 861 *Yersinia* spp., 677 *Serratia*  
208 spp., 186 *Proteus* spp. and 69 *Morganella* spp. genomes (Supplemental Material, Tab. S4). A  
209 less well-conserved homolog of the colibactin determinant was detected in three *F. perrara*  
210 genomes. It should be noted that the number of genomes of *Klebsiella* spp. and *E. coli*  
211 analyzed in this study are markedly higher than those of the other species and lineages due to  
212 the sequencing bias towards *Klebsiella* spp. and *E. coli* strains. Accordingly, a reliable  
213 statement on the prevalence of the colibactin determinant in the different species cannot be  
214 made.

### 215 5.2 Diversity of the colibactin determinant

216 To find out whether the prevalence of the colibactin gene cluster is restricted to specific  
217 phylogenetic lineages of *E. coli* and *Klebsiella* spp., the sequence types of the corresponding *E.*  
218 *coli* and *Klebsiella* spp. isolates were further analyzed. As shown in Fig. S1, the *clb* gene cluster  
219 was enriched in a small subset of *E. coli* STs (twelve out of 11,537 STs, as of 30.10.2020), *K.*  
220 *aerogenes* STs (two out of 214 STs, as of 30.10.2020), and *K. pneumoniae* STs (six out of 5,237  
221 STs, as of 30.10.2020), respectively. In these twelve *E. coli* STs, between 58% and 94% of the  
222 allocated isolates carry the *clb* determinant. A high percentage (ca. 96%) of the *K. aerogenes*  
223 ST4 and ST93 included in our study harbored the colibactin genes. In the tested *K.*  
224 *pneumoniae* strains, all ST3 isolates were *clb*-positive, more than 75% of the analyzed ST23  
225 and ST234 isolates carried the colibactin gene cluster, whereas this was only the case for a  
226 significantly lower percentage of the *K. pneumoniae* isolates allocated to ST11, ST258, and

227 ST48. Table S2 (Supplemental Material) contains a complete list of STs to which colibactin-  
228 positive *E. coli* and *Klebsiella* isolates have been assigned.

229 The nucleotide sequences of the *clb* gene cluster extracted from the 2,169 strains  
230 (Supplemental Material, Tab. S1) were used to generate a recombination-free phylogeny of  
231 the colibactin determinant as shown in Fig. 2A (Fig. S2 and S3 for branch support values and  
232 strain labels/assembly IDs, Fig. S4 indicates predicted recombination events in the *clb* gene  
233 cluster). *Serratia marcescens* strain MSU97 isolated from a plant source, *Erwinia oleae* strain  
234 DAPP-PG531 isolated from olive tree knot, *Klebsiella michiganensis* strain NCTC10261 of an  
235 unknown source, and the *E. coli* phylogroup E strain 14696-7 isolated from the pericardial sac  
236 of a white-tailed deer (*Odocoileus virginianus*) harbor the most genetically distant variants of  
237 the colibactin determinant (Fig. 2A). Within the *Enterobacteriaceae*, a large group of *clb* gene  
238 clusters can be defined, which is dominated by two highly conserved clades present in *E. coli*  
239 phylogenetic lineage B2 and different *Klebsiella* spp. isolates, respectively. The colibactin  
240 determinants detected in *E. cloacae* and *E. hormaechei* belong to the *Klebsiella* clades of *clb*  
241 loci, whereas the *clb* gene clusters found in *C. koseri* and in an unspecified *Salmonella* isolate  
242 represent an independent clade, i.e. *clb6* (Fig. 3). In a few other *E. coli* and *Klebsiella* spp.  
243 genomes, the *clb* determinant can be distinguished from the two major conserved clades of  
244 colibactin determinants observed in *Klebsiella* or *E. coli*. An even more divergent group  
245 comprises the *clb* gene clusters of mainly *E. coli* phylogroup A, B1, and a few B2 isolates, but  
246 also of some *K. pneumoniae* strains (Fig. 3, belonging to *clb* clades *clb1* and *clb2*).

247 Although the *clb* gene clusters found in the large clade of *E. coli* B2 strains are highly  
248 conserved, individual ST-specific lineages, such as *clb10* (from *E. coli* ST141 and ST2015  
249 strains) and *clb12* (from *E. coli* ST121 strains) can still be described within this clade. Beyond  
250 that, we also observed multiple lineages per sequence type, such as *clb2*, *clb11* and *clb16*  
251 found in *E. coli* strains of ST73 (Fig. 3). Three lineages of the *clb* locus were predominantly  
252 detectable in *Klebsiella pneumoniae* strains. They belong to the most distant *Klebsiella*  
253 ST3/ST380 clade (*clb1*), the remaining large and diverse ST258/ST11 clade (*clb8*), and finally  
254 the hypervirulent ST23 clade (*clb9*) (Fig. 3). A phylogeny of the colibactin gene cluster inferred  
255 from concatenated amino acid sequences of the 17 *clb* genes (Fig. S8) was very similar to the  
256 aforementioned recombination-free nucleotide-based phylogeny (Fig. 3).

### 257 **5.3 Diversity of the yersiniabactin determinant in colibactin-positive bacteria**

258 The majority (>98%) of the *clb*-positive enterobacterial strains also harbored the  
259 yersiniabactin genes (*ybt*) (Fig. 3, 3<sup>rd</sup> circle). The *E. coli* strains of phylogroup A, B1, and E as  
260 well as the tested *K. michiganensis* and *E. oleae* strains carrying the most genetically distant  
261 lineages of the colibactin gene cluster, together with the *F. perrara* strains used as an  
262 outgroup, are *ybt*-negative (Fig. 2B). There were no strains that carried multiple copies of the  
263 colibactin gene cluster; yet two well-separated copies of the *ybt* determinant were, however,  
264 found in *C. koseri* strains ATCC BAA-895 and 0123A\_53\_520. It should be noted that the latter  
265 strain is derived from a metagenome (Fig. S6). The phylogenetic analysis indicated that all *ybt*  
266 operons from *E. coli* clustered together (Fig. 2B). Alike the *clb* gene cluster, also the *ybt* locus  
267 of the *E. coli* phylogroup B2 strains was highly conserved. In contrast, the sequence  
268 comparison of the *ybt* determinants of *Klebsiella* spp. resulted in different lineages, which  
269 correlate with lineages *ybt*1, 12, and 17 (ICE*Kp*10), *ybt*9 (ICE*Kp*3), and *ybt*4 (originating from a  
270 plasmid) previously described by Lam and colleagues [38].

### 271 **5.4 Congruent phylogeny of the colibactin and yersiniabactin determinants**

272 The strong coexistence of the colibactin and yersiniabactin determinants on the one hand and  
273 the description of different evolutionary lines of *clb* and *ybt* determinants in different  
274 enterobacterial species on the other hand led us to analyze whether both gene clusters can  
275 predominantly be transferred individually or together. Our results indicate that the clades of  
276 the evolutionary lineages of the *clb* and *ybt* loci are chiefly species/genus-specific. The  
277 phylogeny of *clb* and *ybt* determinants is largely congruent, with the *ybt* gene cluster being,  
278 even more, species/genus-specific than that of the *clb* gene cluster. However, in some strains  
279 we observed evidence of interspecies transfer of these genes: The *clb* and corresponding *ybt*  
280 determinants of the *C. koseri* isolates NCTC10769, ATCC BAA-895 and *E. coli* strains 239A, 926,  
281 GN02370, MOD1-EC5674/5 were allocated to the large *Klebsiella*-dominated lineage *clb*8 (Fig.  
282 3). Additionally, the *clb* gene cluster of *K. pneumoniae* strain k2265 was found within lineage  
283 *clb*6, which predominantly represents *C. koseri* isolates. However, the *ybt* determinant of *K.*  
284 *pneumoniae* strain k2265 belonged to the *ybt*12 lineage represented by *K. oxytoca* isolates.  
285 Similarly, the *ybt* determinant of the aforementioned strains *E. coli* strains GN02370 and *C.*  
286 *koseri* ATCC BAA-895 belonged to lineage *ybt*4 (plasmid originating *ybt* loci) instead of *ybt*17.

287 Regardless of the aforementioned exceptions, clades of the *clb* gene cluster usually correlated  
288 with the corresponding clade of *ybt* genes (Fig. 3, 3<sup>rd</sup> and 4<sup>th</sup> circle).

## 289 5.5 Genetic structure of the MGEs harboring the colibactin determinant

290 To further investigate whether the colibactin and yersiniabactin determinants are jointly  
291 distributed by horizontal gene transfer and to obtain clues to the underlying mechanism, we  
292 compared the chromosomal context of the two polyketide biosynthesis gene clusters and the  
293 genetic structure of associated mobile genetic elements (MGEs). We observed species-specific  
294 structural differences of the chromosomal regions harboring the *clb*, and *ybt* gene clusters  
295 (Fig. 4). In *E. coli* phylogroup B2 strains, the *clb*, and *ybt* determinants are present as part of  
296 two individual pathogenicity islands (PAIs) with their cognate integrase and different tRNA  
297 genes (class I of *clb*-harboring MGE). Both PAIs are located neighboring each other in the  
298 chromosome. Within class Ia, the two PAIs are separated by a type 4 secretion system (T4SS)-  
299 encoding operon (*virB*) and a region that includes two conserved gene sets (Set 1 and Set 2),  
300 different tRNA(*Asn*) loci and an integrase gene (see Table S3 for genes present in the different  
301 conserved gene sets). This region is shown to have been diminished in class Ib structural  
302 variants, where only gene *virB1* of the T4SS determinant is left alongside the conserved  
303 region. In the predominant *E. coli* structural variant, class Ic, the complete T4SS operon  
304 (including *virB1*) has been lost together with gene set 1, the integrase gene was exchanged  
305 and a DNA stretch comprising the genes *yeeO*, tRNA(*Asn*), *cbl* and *gltC* was inverted (shown in  
306 red, Fig. 4). The region separating the *pks* island and the HPI was reduced from a 40-kb (in  
307 class Ia) to a 15-kb stretch in class Ic. In contrast, within all *Klebsiella* strains, the *clb* and *ybt*  
308 gene clusters are part of an ICE, and are separated by a T4SS-encoding operon (*virB*) and the  
309 *hha* gene coding for the hemolysin expression-modulating protein. Downstream of the *ybt*  
310 gene cluster an integrase gene is located followed by a set of genes involved in Fe/Mn/Zn  
311 metabolism (structural class II of *clb*- and *ybt*-harboring chromosomal regions) (Fig. 4).  
312 Interestingly, one type of such ICEs, is located next to genes necessary for microcin E492  
313 biosynthesis (class IIc, Fig. 4). *Enterobacter hormaechei* strains harbor a structurally similar ICE  
314 to that of *Klebsiella* strains. In most *C. koseri* strains, however, the two polyketide  
315 determinants are separated by a large 250-kb chromosomal region. The T4SS-related genes  
316 are closely positioned to the *clb* genes while the gene set involved in Fe/Mn/Zn metabolism is  
317 located downstream of the *ybt* determinant. Only a minor fraction of enterobacterial isolates  
318 analyzed displayed some variation regarding gene content and synteny of these three main

319 classes of colibactin and yersiniabactin-encoding chromosomal regions. The structure of *clb*  
320 and *ybt* regions that do not conform to these major classes are as shown in Figure S5. Instead  
321 of class I, several *E. coli* strains carried class II like chromosomal regions where the T4SS and  
322 the Fe/Mn/Zn metabolism genes were present. In *E. coli* strain HVH128 none of the three  
323 main classes colibactin and yersiniabactin-encoding regions could be identified. Although both  
324 polyketide determinants are co-localized with one integrase gene each, they are widely  
325 separated on this strain's chromosome. *K. pneumoniae* strains TUM14001, TUM14002,  
326 TUM14009, TUM14126, and WCHKP13F2 harbored two T4SS-encoding gene clusters in close  
327 proximity of the *clb* and *ybt* gene clusters and lacked the Fe/Mn/Zn genes, whereas *K.*  
328 *pneumoniae* strain UCI110 was also missing the Fe/Mn/Zn metabolism-related genes. In  
329 contrast to the other *C. koseri* isolates, we detected a class II- instead of a class III-type *clb-ybt*  
330 region in *C. koseri* isolate BAA-895.

## 331 5.6 Organization of the colibactin gene cluster

332 The *clb* gene cluster is composed of 19 genes, which are required for the regulation,  
333 biosynthesis, and transport of colibactin. The origin of this gene cluster is unclear. We,  
334 therefore, compared the structure of the gene clusters representing the homologous *clb* locus  
335 found in *F. perrara* and the phylogenetically most distant and potentially older *clb*  
336 determinants relative to the *E. coli* B2 type of *clb* determinant (Fig. 2A), which are present in *S.*  
337 *marcescens*, *E. oleae*, *K. michiganensis*, and *E. coli* phylogroup E strain 14696-7 (Fig. 5). The  
338 various *clb* determinants correspond in terms of the structure of the genes coding for the  
339 biosynthesis machinery, transport, and resistance to colibactin (*clbB-clbS*) and resemble the  
340 structure of the well-described *clb* locus in *K. pneumoniae/K. aerogenes/E. coli* B2 strains. The  
341 individual *clb* determinants differ more clearly in the presence and localization of the genes  
342 involved in the regulation and activation of colibactin biosynthesis (*clbR* and *clbA*). These two  
343 genes are absent in the homologous gene cluster found in *F. perrara* and in the *clb*  
344 determinant in *S. marcescens*. However, in *F. perrara*, a phosphopantetheinyl transferase  
345 coding for a homolog of ClbA (43% amino acid similarity) and a radical S<sup>1</sup>-adenosylmethionine  
346 (SAM) enzyme-encoding gene are found directly downstream of *clbS*. In *S. marcescens* a helix-  
347 turn-helix (HTH)-type regulatory protein homologous to *clbR* is encoded by a gene located  
348 upstream of *clbB* (78% amino acid similarity). In *E. oleae* and *K. michiganensis* a SAM enzyme-  
349 encoding gene is present directly downstream of *clbA*. Although the colibactin gene clusters in  
350 *E. oleae* and *K. michiganensis* have a high nucleotide sequence similarity of ca. 99.77%, the

351 predicted coding regions of *clbC/D* and *clbH/I/J/K/L/M/N* are noticeably different due to  
352 multiple frameshift deletions in *K. michiganensis*. The structure of the *clb* locus found in  
353 phylogroup E *E. coli* strain 14696-7 already corresponds to the structure of *E. coli* strains of  
354 phylogroup B2 (Fig. 1 and Fig. 5), yet this gene cluster shows the least sequence similarity of  
355 all tested *clb* gene clusters in non-B2 *E. coli* isolates (Fig. 2A) to the determinant occurring in *E.*  
356 *coli* strains of phylogroup B2.

357 Looking at the G+C plot of the colibactin gene clusters, it is obvious that all investigated  
358 enterobacterial *clb* gene clusters show a very similar G+C plot, which has a significantly higher  
359 average G+C content and differs significantly from that of the *clb* homologous gene cluster in  
360 *F. perrara* (Fig. 5). The G+C content profile of these gene clusters indicates that there are two  
361 regions of low G+C content in the enterobacterial *clb* determinants: the region including *clbA*  
362 and *clbR* (at position ca. 1,500 bp - 3,000 bp of the colibactin gene cluster) and the region  
363 spanning *clbD* and *clbE* (at position ca. 15,000 bp – 16,500 bp of the *clb* gene cluster). The G+C  
364 content drop in the region including *clbA* and *clbR* (at position ca. 1,500 bp - 3,000 bp of the  
365 colibactin gene cluster) is associated with a predicted recombination site, which is located  
366 upstream of or interrupting *clbB* (Fig. S4).

367 The comparison of structural features of the *clb* gene cluster also included the VNTR region  
368 located upstream of *clbR* in the *clbR-clbB* intergenic region. The size of the VNTR region has  
369 been described to range from 2 to 20 (Putze *et al.*, 2009). The VNTR copy number distribution  
370 in ca. 1,300 *clb*-positive genomes demonstrated that there is a preference for VNTR regions  
371 ranging from 7-10 copy numbers. Copy numbers from 18-34 were present in only a few strains  
372 (Figure S6). Species and/or ST-specific copy number variation was not observed.

373 Comparative genomic analysis of multiple colibactin-encoding determinants based on (draft)  
374 genome sequences led to the observation that the homologous genes *clbJ* and *clbK* are prone  
375 to fusion/deletion (Lam *et al.*, 2018). We also observed that in several assemblies of the *clb*  
376 gene cluster 625 bp from the 3' end of *clbJ* and 3,540 bp from the 5' end of *clbK* including the  
377 11 bp intergenic region are missing, for a total of 4,174 bp (Fig. S5). The assembly of our  
378 internally generated genome sequences produced by short read (Illumina) sequencing  
379 showed this *clbJ/K* fusion/deletion. However, since assemblies of sequence data of the same  
380 strains generated by a long-read sequencing technology (PacBio), where the long reads  
381 covered both genes, had both *clbJ* and *clbK* completely present, we assume that the *clbJ/K*

382 fusions described are artificial and result from erroneous assemblies of short-read sequencing  
383 data.

## 384 5.7 Quantification of colibactin synthesis in selected strains

385 To investigate a possible correlation between the genetic structure of the *clb* determinant or  
386 the genetic background of the corresponding host strain with colibactin expression, we  
387 quantified N-myristoyl-D-asparagine levels produced *in vitro* by selected *clb*-positive *E. coli*,  
388 *Klebsiella* spp. and *C. koseri* strains covering the diversity of *clb* determinants in these species  
389 (Fig. 6). Based on the detected relative amount of N-myristoyl-D-asparagine produced, the  
390 investigated isolates can be roughly divided into two groups: One group included most of the  
391 measured *E. coli* strains that produced only very low relative amounts of N-myristoyl-D-  
392 asparagine. In contrast, the tested *C. koseri*, *K. aerogenes*, and *K. pneumoniae* isolates and the  
393 *E. coli* isolates CFT073 and N1 showed a 3 to 70-fold higher N-myristoyl-D-asparagine  
394 production. Also within the *C. koseri*, *K. aerogenes*, and *K. pneumoniae* isolates, we found  
395 differences in the relative N-myristoyl-D-asparagine levels. However, these differences were  
396 not as strong as among the eight *E. coli* isolates studied. The observed relative N-myristoyl-D-  
397 asparagine levels do not indicate phylogroup, ST, or species-specific differences in colibactin  
398 production. For example, the *E. coli* strain 1873, although the *clb* gene cluster present in this  
399 strain is phylogenetically more closely related to that of *E. coli* strain N1, shows a significantly  
400 weaker N-myristoyl-D-asparagine production than *E. coli* N1. Similarly, it should be noted that  
401 *C. koseri* MFP3 produces less N-myristoyl-D-asparagine than other closely related *C. koseri*.

402

## 403 6 Discussion

### 404 *Prevalence and mobility of the colibactin gene cluster in Enterobacteriaceae*

405 The ability of MGEs to be exchanged between and within species plays a major role in the  
406 extent and speed of microbial evolution. Because MGEs are known to emerge and evolve  
407 separately from the host, it is important to explain the development and diversity of MGEs  
408 independently. In our study, we investigated on the one hand the nucleotide sequence  
409 variability of the colibactin gene cluster with the associated yersiniabactin determinants, and  
410 on the other hand the structural diversity of the MGEs hosting the colibactin and  
411 yersiniabactin determinants responsible for their horizontal distribution. Both polyketide  
412 determinants together could be detected in certain members of the *Enterobacteriaceae*,

413 mainly in *E. coli*, *K. pneumoniae*, *K. aerogenes*, and *C. koseri*, but also in *Enterobacter cloacae*,  
414 *Enterobacter hormachei*, and possibly in uncharacterized *Salmonella* sp. isolates. The  
415 colibactin, but not the yersiniabactin determinants could also be detected in *K. michiganensis*,  
416 *E. oleae*, and *S. marcescens* and a few phylogroup A and B1 *E. coli* isolates (Fig. 2A).  
417 Interestingly, the colibactin genes have not yet been described in not only archaeal genomes  
418 but also *S. enterica* and *S. bongori* genomes, although we have screened more than 112,000  
419 *Salmonella* spp. genomes. Overall, it is remarkable that the *clb* gene cluster was only found in  
420 extraintestinal clinical or in fecal isolates of healthy hosts, but not in enterobacterial  
421 diarrhoeal pathogens such as *Yersinia* spp., *Salmonella* spp., *Shigella* spp. as well as the  
422 various intestinal pathogenic *E. coli* pathotypes. This observation is consistent with previously  
423 published data [2, 50]. In this context, it would be interesting to study in the future whether  
424 and in which way colibactin expression supports extraintestinal pathogenicity or intestinal  
425 persistence and colonization, but is detrimental to the pathogenesis of diarrhoeal pathogens.  
426 Possibly, the genetic background also plays an important role in the horizontal distribution  
427 and establishment of MGEs carrying both polyketide determinants. The fact that the  
428 colibactin determinant has so far been preferentially distributed in only some, often highly  
429 virulent STs in *E. coli* and *K. pneumoniae* or *K. aerogenes* and not more broadly within the  
430 respective species (Fig. 2) [38, 39], could also indicate that the transmission, uptake or  
431 chromosomal integration of these MGEs is restricted. It is interesting to note that the *clb* gene  
432 cluster as a whole is highly conserved and usually characteristic of the respective species or  
433 genus (Fig. 3). Nevertheless, ST-specific variants have been found within a species, e.g. in *K.*  
434 *pneumoniae* ST3 and ST23. Several groups of sequence variants of the *clb* gene cluster can  
435 also be found within one ST, such as in the *E. coli* ST73 and ST95 (Fig. 3). These results show  
436 that, on the one hand, intraspecies transfer of the colibactin determinant can happen, but on  
437 the other hand, certain adaptations of the *clb* genes to a specific genetic background can also  
438 occur at the nucleotide level. In addition, examples of interspecies transfer of the *clb* and *ybt*  
439 genes can be seen between *Klebsiella* spp. and *E. hormachei* strains. Similarly and in contrast  
440 to the majority of *E. coli* isolates, we found *clb* and *ybt* gene clusters in some *E. coli* isolates,  
441 which are assigned to the large clade of *Klebsiella/Enterobacter/Citrobacter*-specific *clb* and  
442 *ybt* variants (Fig. 2A and 3). Apart from the interspecies transfer of the entire *clb*- and *ybt*-  
443 containing MGE, we also identified an example that shows that the two polyketide gene  
444 clusters can also be exchanged independently, as in the case of the *clb* gene cluster of *K.*



445 *pneumoniae* strain k2265, which belongs to colibactin clade *clb6* (predominantly *C. koseri*  
446 lineage of *clb* loci), whereas the *ybt* determinant in this strain is assigned to the yersiniabactin  
447 clade *ybt12* instead of *ybt3*, which is usually associated with *C. koseri* strains carrying clade  
448 *clb6* (Fig. 3).

#### 449 *Structural diversity of the colibactin-yersiniabactin region*

450 The structural analysis of the genomic region comprising the *clb* and *ybt* determinants in their  
451 chromosomal sequence context is an important aspect to understand the evolution of these  
452 polyketide determinants and their origin. In principle, three structural constellations (classes I  
453 to III) can be described, in which the *clb* and *ybt* gene clusters are present (Fig. 4). Class I  
454 depicts the *clb* and *ybt* gene clusters found in the majority of *E. coli* isolates, each associated  
455 with a tRNA(*Asn*) and an integrase gene. In class Ia, the *clb* genes are chromosomally inserted  
456 at the tRNA locus *asnV*. Our analyses suggest that in class I, the type 4 secretion system gene  
457 cluster ("mobilization module") and conserved neighboring genes (set 1) have been lost in a  
458 stepwise process, from class Ia to Ic. A further structural modification in this region is  
459 represented by the inversion of the *yeeO*-tRNA(*Asn*)-*cbl-gltC* gene set (Fig. 4, red arrow), as a  
460 result of which in class Ic, the tRNA gene *asnW* is located closest to the *clb* genes. Taking into  
461 account that the *E. coli* strains with class Ia and Ib structures are found in the potentially  
462 earlier phylogenetic clades, we hypothesize that the MGE harboring the *clb* genes was  
463 introduced into the *E. coli* chromosome separately from that carrying the *ybt* determinant.  
464 Both MGEs were then progressively modified as described above. We do not yet have an  
465 explanation why the resulting class Ic colibactin-yersiniabactin region, which has been  
466 described as two pathogenicity islands (PAIs) comprising the colibactin and yersiniabactin  
467 determinants, respectively [1, 2], is only found in phylogroup B2 strains and only there has it  
468 become so successful.

469 Class II includes different variants of an ICE, in which the two polyketide determinants are  
470 present in association with a type 4 secretion system-encoding "mobilization module" and a  
471 "module" consisting of genes that contribute to Fe/Mn/Zn metabolism (Fig. 4, Supplementary  
472 Table S3). This type of ICE was found in *Klebsiella* spp. and *E. hormachei* and in a few cases in  
473 *E. coli* and *C. koseri*. Unlike in class I and III, the ICE in class II does not only have different  
474 tRNA(*Asn*) loci serving as chromosomal insertion sites, but also lacks tRNA(*Asn*) and integrase  
475 genes in between the *clb* and *ybt* genes. In a population-wide analysis of *Klebsiella* spp.  
476 strains, this ICE was designated ICE*Kp10* and described as being associated with different

477 combinations of *ybt* and *clb* gene lineages [38]. In contrast to class I and II colibactin-  
478 yersiniabactin regions, the *clb* and *ybt* gene clusters are located far apart on the chromosome  
479 in most of the *C. koseri* genomes studied (class III) (Fig. 4).

480 The existence of an ICE that unites the *clb* and *ybt* gene clusters is the easiest way to explain  
481 the co-localization and the joint transfer of the two determinants and thus the high  
482 correlation of *clb* and *ybt* phylogenetic clades (Fig. 3, 4th circle) in *Klebsiella* spp. strains.  
483 Despite slight differences in the sequence context and different chromosomal insertion sites  
484 (Fig. 4), the ICEs of the four class II variants have an overall identical genetic structure (Fig. 4).  
485 The uptake of this ICE thus leads to the acquisition of both the *clb* and the *ybt* gene clusters.  
486 The presence of the Fe/Mn/Zn metabolic genes neighboring these ICE variants with an  
487 additional integrase gene indicates recombination processes that can alter the genetic  
488 structure of the ICE. The clear separation between the *clb*/T4SS module and the *ybt*-Fe/Mn/Zn  
489 metabolism module in the *C. koseri* genomes points towards rearrangement/relocation of the  
490 *ybt* region, after ICE integration into the chromosome. The fact that *C. koseri* strain ATCC BAA-  
491 895 possesses in addition to the complete ICE (class II) a second *ybt* gene cluster (99.98%  
492 nucleotide similarity to the *ybt* genes present in the ICE) that is located far away from the ICE  
493 (Fig. S6), supports the hypothesis that the individual polyketide gene clusters can also be  
494 integrated into the genome independently of each other. This state could result, for example,  
495 from the initial chromosomal integration of different ICE variants, as described by Lam and  
496 colleagues in *Klebsiella* spp. [38]. As a result of deletion events, through which individual  
497 modules are subsequently deleted from one of the two ICEs, the second copy of the *ybt* gene  
498 cluster remains in the genome as a fragment of the degenerated ICE. The presence of two  
499 non-identical T4SS modules in the *K. pneumoniae* strains TUM14001, TUM14002, TUM14009,  
500 TUM14126, and WCHKP13F2 (Fig. S6 and Fig. 3, *clb4*) associated with the *clb* or *ybt* module  
501 (Fig. S6) could be the result of such degeneration of different ICEs. In this way, our  
502 observations on the phylogeny (Fig. 2) and structure (Fig. 4) of the two polyketide  
503 determinants and their sequence context can be reconciled, which show that despite the  
504 predominant species/genus specificity, there are also sequence type-specific lineages of the  
505 *clb* genes, which do not necessarily have to match that of the associated *ybt* genes.

506 In this context, one could imagine that the arrangement of *clb* and *ybt* determinants in *E. coli*  
507 strains (class I) also results from independent integration events of different MGEs, which  
508 subsequently degenerated as a result of a stabilization process of these MGEs [26, 28]. This

509 premise is supported by the absence of *ybt* genes in the *clb*-positive *E. coli* strains, which carry  
510 the phylogenetically most distant *clb* determinants compared to the *clb* genes of phylogroup  
511 B2 isolates (Fig. 2A), along with the presence of integrase and tRNA(*Asn*) genes in close  
512 proximity to both polyketide determinants (Fig. 4, class I).

### 513 *Evolution of the colibactin determinant*

514 Homologs of the *clb* gene cluster were detected in marine alpha-proteobacteria such as  
515 various *Pseudovibrio* sp. (isolates AD26, FO-BEG1, POLY-S9) and *Pseudovibrio denitrificans*  
516 (isolates DSM 17465 and JCM12308) [51]. Despite the general conformity of the genetic  
517 structure of these gene clusters, their nucleotide sequence identity to the colibactin gene  
518 cluster is quite low (<26%). Therefore, it was hypothesized that these *Pseudovibrio* isolates  
519 have the potential to produce molecules related to colibactin [51]. Another homolog of the  
520 colibactin determinant with a higher (62%) amino acid sequence identity is found in *F. perrara*  
521 [24]. While the genes required for biosynthesis and transport of the polyketide are present,  
522 the genes corresponding to *clbA* and *clbR* are missing in this gene cluster (Fig. 5). The case is  
523 similar with the colibactin gene cluster in *S. marcescens*. While in *F. perrara*, a gene coding for a  
524 *clbA* homolog and a gene coding for a SAM enzyme are located immediately downstream of  
525 the *clbS* homolog, a *clbR* homolog is located upstream of *clbB* in *S. marcescens*. It is therefore  
526 conceivable that the *clbA* homolog in *F. perrara* and the *clbR* homolog in *S. marcescens* are  
527 involved in the activation or regulation of colibactin biosynthesis in these bacteria (Fig. 5). It  
528 has been described that S-adenosylmethionine (SAM) is used in NRPS modules for colibactin  
529 biosynthesis [52]. Looking at the genetic structure of the *clbB-S* homologous genes cluster in  
530 *F. perrara* and the *clb* gene cluster in *S. marcescens*, one can assume that genes involved in  
531 the regulation and activation of the biosynthetic pathway including the SAM enzyme gene as  
532 well as *clbA* and *clbR* homologs have been fused upstream to the already existing part of the  
533 island (*clbB-clbS*) to improve regulation of polyketide biosynthesis. Without having further  
534 knowledge about the origin of the colibactin biosynthetic genes themselves, the acquisition of  
535 the regulatory/activating genes is obviously among the last evolutionary steps that led to the  
536 structural organization of what we currently describe as the colibactin gene cluster. This  
537 hypothesis is supported by the abrupt decrease in G+C content and the presence of the  
538 predicted recombination (Fig. S4) directly upstream of *clbB* in many *clb* determinants.  
539 Furthermore, a module consisting of a gene for a SAM enzyme and a *clbA* homolog is not only  
540 located directly downstream from the *clbS* homolog in *F. perrara*, but also upstream from *clbR*

541 in *K. michiganensis* and *E. oleae*, which represent evolutionarily older variants of colibactin-  
542 positive *Enterobacteriaceae* (Fig. 5).

#### 543 *Expression of colibactin in different hosts*

544 Furthermore, we investigated the question of how differently colibactin is expressed within  
545 enterobacterial genera or even within different lineages of the same species. Interestingly, we  
546 observed an often lower production level of N-myristoyl-D-asparagine in *E. coli* isolates  
547 compared to *K. aerogenes*, *K. pneumoniae*, and *C. koseri* (Fig. 6), which may be expected since  
548 *E. coli* is described as a non-optimal producer of complex secondary metabolites [53].  
549 However, it is of interest that the amount of N-myristoyl-D-asparagine produced in *E. coli*  
550 strains CFT073 and N1 is comparable to that of other enterobacterial genera (Fig. 6). A  
551 species- or lineage-specific ability to produce N-myristoyl-D-asparagine could not be  
552 determined so far. Future studies will have to investigate which bacterial factors are  
553 important for colibactin production and how the strain-specific differences in the expression  
554 of this polyketide come about. The systematic comparison of phenotypic colibactin  
555 production with information on the genomic context, regulatory and metabolic properties of  
556 host strains, and their classification in a phylogenetic context should help us to identify  
557 bacterial factors that affect colibactin synthesis.

## 558 **7 Conclusion**

559 The colibactin and yersiniabactin gene clusters are highly conserved polyketide determinants  
560 within some *Enterobacteriaceae*. They usually coexist together in the genome and are also  
561 linked to each other at the biosynthetic level. With the exception of *E. coli*, the two gene  
562 clusters are part of an ICE, which allows the horizontal transfer of both secondary metabolite  
563 determinants usually within one species/genus. Bacteria of the genus *Klebsiella* played an  
564 important role in the evolution and distribution of both gene clusters. A large number of  
565 different ICEs has been described in *Klebsiella spp.*, which besides several other groups of  
566 genes include the yersiniabactin determinant [38]. Recombination and rearrangements events  
567 between different ICE types may have contributed to the evolution of the ICE variants so far  
568 identified in *Klebsiella spp.* and other enterobacteria as well as to the further degeneration of  
569 such MGEs leading to the colibactin and yersiniabactin-encoding PAIs present in phylogroup  
570 B2 *E. coli* strains (Fig. 7).

571 The phylogeny of the *clb* determinants does not determine the level of phenotypic colibactin  
572 production. The underlying bacterial factors responsible for the colibactin production  
573 efficiency of individual strains need to be identified in future work.

574 Our investigations provide deeper insights into the evolution of the colibactin gene cluster in  
575 *Enterobacteriaceae*. Based on our findings, we can extend the current explanation for the co-  
576 existence and genetic co-localization of both gene clusters. The combination of a PPTase-  
577 encoding gene (*clbA*) with the *clbB-S* biosynthetic gene cluster during the evolution of the *clb*  
578 determinant not only enabled the efficient activation of the colibactin biosynthesis machinery,  
579 but also linked the colibactin and yersiniabactin determinants, which are functionally  
580 connected by the activity of PPTase ClbA. This enables the bacteria to synthesize both  
581 functionally different secondary metabolites, which leads to a stabilization of the co-existence  
582 and co-localization of the two gene clusters in the genome. Our data underpin the importance  
583 of mobile genetic elements, especially of ICEs, for genomic diversity and variability in  
584 enterobacteria as well as for the evolution of more complex bacterial phenotypes, such as the  
585 combined expression of the secondary metabolites colibactin and yersiniabactin.

## 586 **8 Author statements**

### 587 **8.1 Authors and contributors**

588 H.W., A.W. and U.D. conceptualized the project. H.W., A.W., D.S. ran the analyses. R.M., M.S., and E.O.  
589 contributed reagents and new tools. H.W., A.W., and U.D. analysed the data. H.W., A.W. and U.D.  
590 wrote the manuscript. H.W., A.W., M.S., R.B., E.O., R.M. and U.D. edited and revised the manuscript.  
591 All authors read, commented on and approved the final manuscript.

### 592 **8.2 Conflicts of interest**

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

### 594 **8.3 Funding information**

595 The work of the Münster team was supported by the German Research Foundation (grant  
596 DO789/11-1 and grant 281125614/GRK2220 (EvoPAD project A3)).

### 597 **8.4 Acknowledgments**

598 Data reported in this study appear in part in the PhD theses of A. Wallenstein and H. Wami.  
599 We thank M.K. Mammel (United States Food and Drug Administration) for providing *E. coli*

600 strain N1. We thank K. Tegelkamp and O. Mantel (Münster) for excellent technical support.  
601 We gratefully acknowledge CPU time on the high performance cluster PALMA@WWU  
602 Münster.

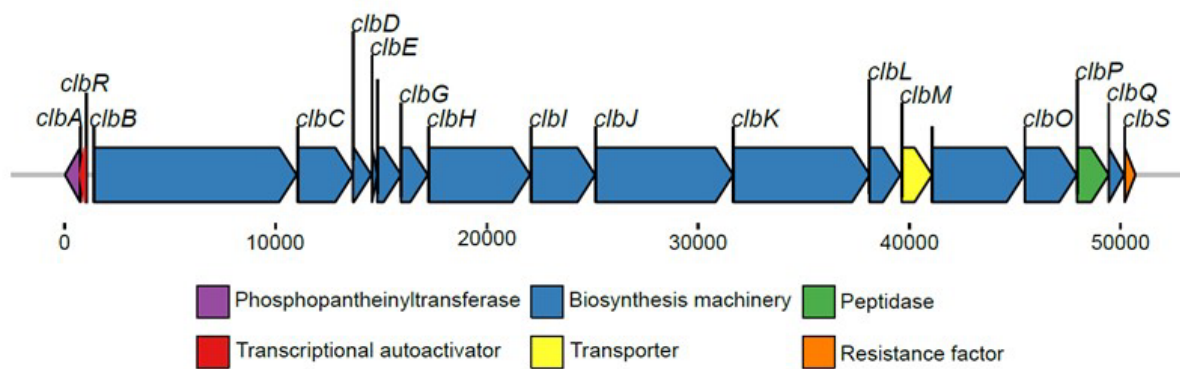
## 603 9 References

- 604 1. **Nougayrède JP, Homburg S, Taieb F, Boury M, Brzuszkiewicz E et al.** *Escherichia coli* induces  
605 DNA double-strand breaks in eukaryotic cells. *Science* 2006;313(5788):848-851.
- 606 2. **Putze J, Hennequin C, Nougayrède JP, Zhang W, Homburg S et al.** Genetic structure and  
607 distribution of the colibactin genomic island among members of the family *Enterobacteriaceae*.  
608 *Infect Immun* 2009;77(11):4696-4703.
- 609 3. **Xue M, Kim CS, Healy AR, Wernke KM, Wang Z et al.** Structure elucidation of colibactin and its  
610 DNA cross-links. *Science*, 2019;365(6457).
- 611 4. **Homburg S, Oswald E, Hacker J, Dobrindt U.** Expression analysis of the colibactin gene cluster  
612 coding for a novel polyketide in *Escherichia coli*. *FEMS Microbiol Lett* 2007;275(2):255-262.
- 613 5. **Martin P, Marcq I, Magistro G, Penary M, Garcie C et al.** Interplay between siderophores and  
614 colibactin genotoxin biosynthetic pathways in *Escherichia coli*. *PLoS Pathog* 2013;9(7):e1003437.
- 615 6. **Wallenstein A, Rehm N, Brinkmann M, Selle M, Bossuet-Greif N et al.** ClbR Is the Key  
616 Transcriptional Activator of Colibactin Gene Expression in *Escherichia coli*. *mSphere*  
617 2020;5(4):e00591-00520.
- 618 7. **Healy AR, Wernke KM, Kim CS, Lees NR, Crawford JM et al.** Synthesis and reactivity of  
619 precolibactin 886. *Nat Chem* 2019;11(10):890-898.
- 620 8. **Li ZR, Li J, Cai W, Lai JYH, McKinnie SMK et al.** Macrocyclic colibactin induces DNA double-strand  
621 breaks via copper-mediated oxidative cleavage. *Nat Chem* 2019;11(10):880-889.
- 622 9. **Li ZR, Li J, Gu JP, Lai JY, Duggan BM et al.** Divergent biosynthesis yields a cytotoxic  
623 aminomalonate-containing precolibactin. *Nat Chem Biol* 2016;12(10):773-775.
- 624 10. **Bossuet-Greif N, Vignard J, Taieb F, Mirey G, Dubois D et al.** The colibactin genotoxin generates  
625 DNA interstrand cross-links in infected cells. *mBio* 2018;9(2):e02393-02317.
- 626 11. **Cuevas-Ramos G, Petit CR, Marcq I, Boury M, Oswald E et al.** *Escherichia coli* induces DNA  
627 damage *in vivo* and triggers genomic instability in mammalian cells. *Proc Natl Acad Sci USA*  
628 2010;107(25):11537-11542.
- 629 12. **Buc E, Dubois D, Sauvanet P, Raisch J, Delmas J et al.** High prevalence of mucosa-associated *E.*  
630 *coli* producing cyclomodulin and genotoxin in colon cancer. *PLoS One*, 2013;8(2):e56964.
- 631 13. **Cougnoux A, Dalmaso G, Martinez R, Buc E, Delmas J et al.** Bacterial genotoxin colibactin  
632 promotes colon tumour growth by inducing a senescence-associated secretory phenotype. *Gut*  
633 2014;63(12):1932-1942.
- 634 14. **Dalmaso G, Cougnoux A, Delmas J, Darfeuille-Michaud A, Bonnet R.** The bacterial genotoxin  
635 colibactin promotes colon tumor growth by modifying the tumor microenvironment. *Gut*  
636 *Microbes* 2014;5(5):675-680.
- 637 15. **Dziubanska-Kusibab PJ, Berger H, Battistini F, Bouwman BAM, Iftekhar A et al.** Colibactin DNA-  
638 damage signature indicates mutational impact in colorectal cancer. *Nat Med* 2020;26(7):1063-  
639 1069.
- 640 16. **Marcq I, Martin P, Payros D, Cuevas-Ramos G, Boury M et al.** The genotoxin colibactin  
641 exacerbates lymphopenia and decreases survival rate in mice infected with septicemic  
642 *Escherichia coli*. *J Infect Dis*, 2014;210(2):285-294.
- 643 17. **Pleguezuelos-Manzano C, Puschhof J, Rosendahl Huber A, van Hoeck A, Wood HM et al.**  
644 Mutational signature in colorectal cancer caused by genotoxic pks(+) *E. coli*. *Nature*  
645 2020;580(7802):269-273.
- 646 18. **Faïs T, Delmas J, Barnich N, Bonnet R, Dalmaso G.** Colibactin: More Than a New Bacterial  
647 Toxin. *Toxins* 2018;10(4).
- 648 19. **Wassenaar TM.** *E. coli* and colorectal cancer: a complex relationship that deserves a critical  
649 mindset. *Crit Rev Microbiol*, Review 2018;44(5):619-632.
- 650 20. **Wilson MR, Jiang Y, Villalta PW, Stornetta A, Boudreau PD et al.** The human gut bacterial  
651 genotoxin colibactin alkylates DNA. *Science*, 2019;363(6428).

- 652 21. **Brzuszkiewicz E, Brüggemann H, Liesegang H, Emmerth M, Olschlager T et al.** How to become a  
653 uropathogen: comparative genomic analysis of extraintestinal pathogenic *Escherichia coli*  
654 strains. *Proc Natl Acad Sci U S A* 2006;103(34):12879-12884.
- 655 22. **Bellanger X, Payot S, Leblond-Bourget N, Guedon G.** Conjugative and mobilizable genomic  
656 islands in bacteria: evolution and diversity. *FEMS Microbiol Rev* 2014;38(4):720-760.
- 657 23. **Schubert S, Rakin A, Karch H, Carniel E, Heesemann J.** Prevalence of the "high-pathogenicity  
658 island" of *Yersinia* species among *Escherichia coli* strains that are pathogenic to humans. *Infect*  
659 *Immun* 1998;66(2):480-485.
- 660 24. **Engel P, Vizcaino MI, Crawford JM.** Gut symbionts from distinct hosts exhibit genotoxic activity  
661 via divergent colibactin biosynthesis pathways. *Appl Environ Microbiol* 2015;81(4):1502-1512.
- 662 25. **Moretti C, Hosni T, Vandemeulebroecke K, Brady C, De Vos P et al.** *Erwinia oleae* sp. nov.,  
663 isolated from olive knots caused by *Pseudomonas savastanoi* pv. *savastanoi*. *Int J Syst Evol*  
664 *Microbiol* 2011;61(Pt 11):2745-2752.
- 665 26. **Hacker J, Kaper JB.** Pathogenicity islands and the evolution of microbes. *Annu Rev Microbiol*  
666 2000;54:641-679.
- 667 27. **Messerer M, Fischer W, Schubert S.** Investigation of horizontal gene transfer of pathogenicity  
668 islands in *Escherichia coli* using next-generation sequencing. *PLoS One* 2017;12(7):e0179880.
- 669 28. **Schneider G, Dobrindt U, Middendorf B, Hochhut B, Szjarto V et al.** Mobilisation and  
670 remobilisation of a large archetypal pathogenicity island of uropathogenic *Escherichia coli* in  
671 *vitro* support the role of conjugation for horizontal transfer of genomic islands. *BMC Microbiol*  
672 2011;11:210.
- 673 29. **Chen YT, Lai YC, Tan MC, Hsieh LY, Wang JT et al.** Prevalence and characteristics of *pks*  
674 genotoxin gene cluster-positive clinical *Klebsiella pneumoniae* isolates in Taiwan. *Sci Rep*  
675 2017;7:43120.
- 676 30. **Dubois D, Delmas J, Cady A, Robin F, Sivignon A et al.** Cyclomodulins in urosepsis strains of  
677 *Escherichia coli*. *J Clin Microbiol* 2010;48(6):2122-2129.
- 678 31. **Johnson JR, Johnston B, Kuskowski MA, Nougayrede JP, Oswald E.** Molecular epidemiology and  
679 phylogenetic distribution of the *Escherichia coli pks* genomic island. *J Clin Microbiol*  
680 2008;46(12):3906-3911.
- 681 32. **Lai YC, Lin AC, Chiang MK, Dai YH, Hsu CC et al.** Genotoxic *Klebsiella pneumoniae* in Taiwan.  
682 *PLoS One* 2014;9(5):e96292.
- 683 33. **Micenikova L, Benova A, Frankovicova L, Bosak J, Vrba M et al.** Human *Escherichia coli* isolates  
684 from hemocultures: Septicemia linked to urogenital tract infections is caused by isolates  
685 harboring more virulence genes than bacteraemia linked to other conditions. *Int J Med*  
686 *Microbiol* 2017;307(3):182-189.
- 687 34. **Nowrouzian FL, Oswald E.** *Escherichia coli* strains with the capacity for long-term persistence in  
688 the bowel microbiota carry the potentially genotoxic *pks* island. *Microb Pathog*, 2012;53(3-  
689 4):180-182.
- 690 35. **Krieger JN, Dobrindt U, Riley DE, Oswald E.** Acute *Escherichia coli* prostatitis in previously health  
691 young men: bacterial virulence factors, antimicrobial resistance, and clinical outcomes. *Urology*,  
692 2011;77(6):1420-1425.
- 693 36. **McCarthy AJ, Martin P, Cloup E, Stabler RA, Oswald E et al.** The genotoxin colibactin is a  
694 determinant of virulence in *Escherichia coli* K1 experimental neonatal systemic infection. *Infect*  
695 *Immun*, 2015;83(9):3704-3711.
- 696 37. **Yoshikawa Y, Tsunematsu Y, Matsuzaki N, Hirayama Y, Higashiguchi F et al.** Characterization of  
697 colibactin-producing *Escherichia coli* isolated from Japanese patients with colorectal cancer. *Jpn*  
698 *J Infect Dis* 2020.
- 699 38. **Lam MMC, Wick RR, Wyres KL, Gorrie CL, Judd LM et al.** Genetic diversity, mobilisation and  
700 spread of the yersiniabactin-encoding mobile element ICEKp in *Klebsiella pneumoniae*  
701 populations. *Microb Genom* 2018;4(9).



- 702 39. **Shen P, Berglund B, Chen Y, Zhou Y, Xiao T et al.** Hypervirulence Markers Among Non-ST11  
703 Strains of Carbapenem- and Multidrug-Resistant *Klebsiella pneumoniae* Isolated From Patients  
704 With Bloodstream Infections. *Front Microbiol* 2020;11:1199.
- 705 40. **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M et al.** SPAdes: a new genome  
706 assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* 2012;19(5):455-  
707 477.
- 708 41. **Seemann T.** Prokka: rapid prokaryotic genome annotation. *Bioinformatics*, 2014;30(14):2068-  
709 2069.
- 710 42. **Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J et al.** BLAST+: architecture and  
711 applications. *BMC Bioinformatics* 2009;10:421.
- 712 43. **Blin K, Pascal Andreu V, de Los Santos ELC, Del Carratore F, Lee SY et al.** The antiSMASH  
713 database version 2: a comprehensive resource on secondary metabolite biosynthetic gene  
714 clusters. *Nucleic Acids Res* 2019;47(D1):D625-D630.
- 715 44. **Assefa S, Keane TM, Otto TD, Newbold C, Berriman M.** ABACAS: algorithm-based automatic  
716 contiguation of assembled sequences. *Bioinformatics* 2009;25(15):1968-1969.
- 717 45. **Lassmann T, Sonnhammer EL.** Kalign - an accurate and fast multiple sequence alignment  
718 algorithm. *BMC Bioinformatics* 2005;6:298.
- 719 46. **Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA et al.** Rapid phylogenetic analysis of  
720 large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids*  
721 *Res* 2015;43(3):e15.
- 722 47. **Stamatakis A.** RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
723 phylogenies. *Bioinformatics* 2014;30(9):1312-1313.
- 724 48. **Benson G.** Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*  
725 1999;27(2):573-580.
- 726 49. **Bian X, Fu J, Plaza A, Herrmann J, Pistorius D et al.** *In vivo* evidence for a prodrug activation  
727 mechanism during colibactin maturation. *Chembiochem* 2013;14(10):1194-1197.
- 728 50. **Morgan RN, Saleh SE, Farrag HA, Aboulwafa MM.** Prevalence and pathologic effects of  
729 colibactin and cytotoxic necrotizing factor-1 (Cnf 1) in *Escherichia coli*: experimental and  
730 bioinformatics analyses. *Gut Pathog* 2019;11:22.
- 731 51. **Naughton LM, Romano S, O'Gara F, Dobson ADW.** Identification of Secondary Metabolite Gene  
732 Clusters in the *Pseudovibrio* Genus Reveals Encouraging Biosynthetic Potential toward the  
733 Production of Novel Bioactive Compounds. *Front Microbiol* 2017;8:1494.
- 734 52. **Zha L, Jiang Y, Henke MT, Wilson MR, Wang JX et al.** Colibactin assembly line enzymes use S-  
735 adenosylmethionine to build a cyclopropane ring. *Nat Chem Biol* 2017;13(10):1063-1065.
- 736 53. **Zhang H, Wang Y, Pfeifer BA.** Bacterial hosts for natural product production. *Mol Pharm*  
737 2008;5(2):212-225.
- 738
- 739

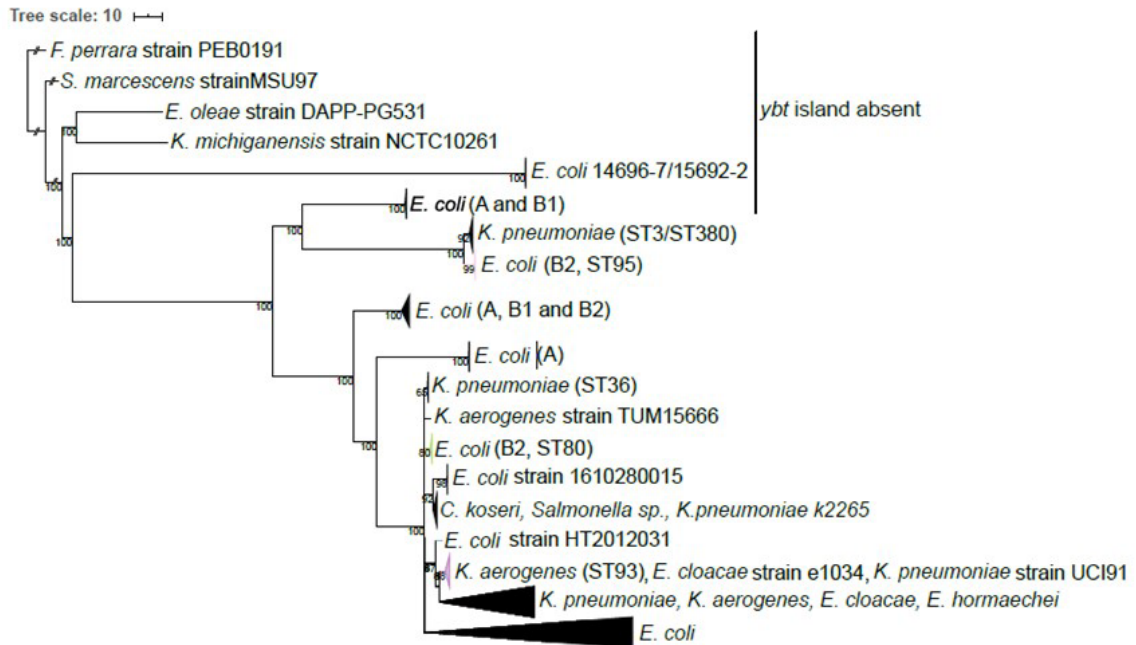


740

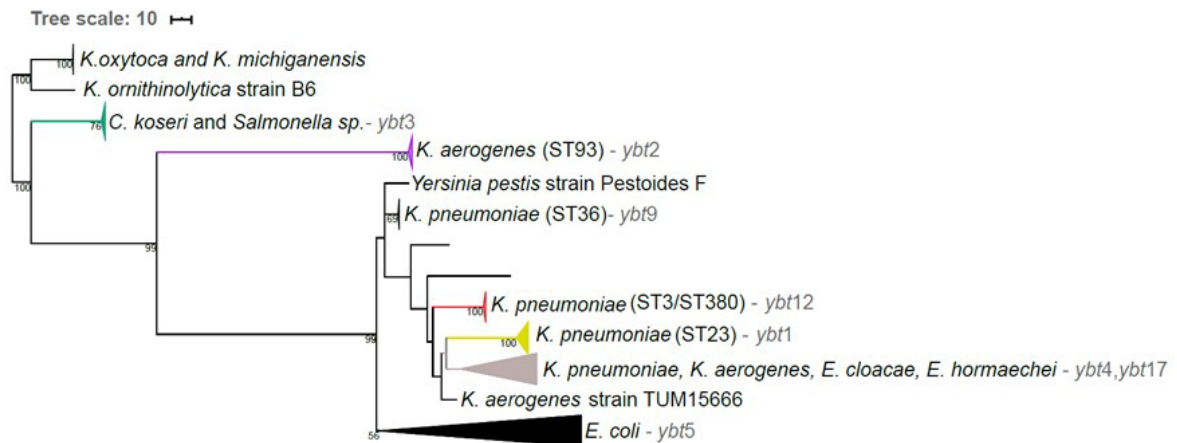
741 **Figure 1.** Schematic representation of the genomic architecture of the *pks* island (ca. 54 kb)  
742 present in *E. coli* strain M1/5. The 19 genes within the island are colored with respect to their  
743 function. The island codes for a phosphopantetheinyl transferase (*clbA*, in purple), a  
744 transcriptional autoactivator (*clbR*, in red), multiple core biosynthetic genes (*clbB*-*clbL*, *clbN*,  
745 *clbO*, and *clbQ* in blue), a transporter (*clbM*, yellow), a peptidase (*clbP*, in green) and a  
746 resistance factor (*clbS*, orange).

747

A)

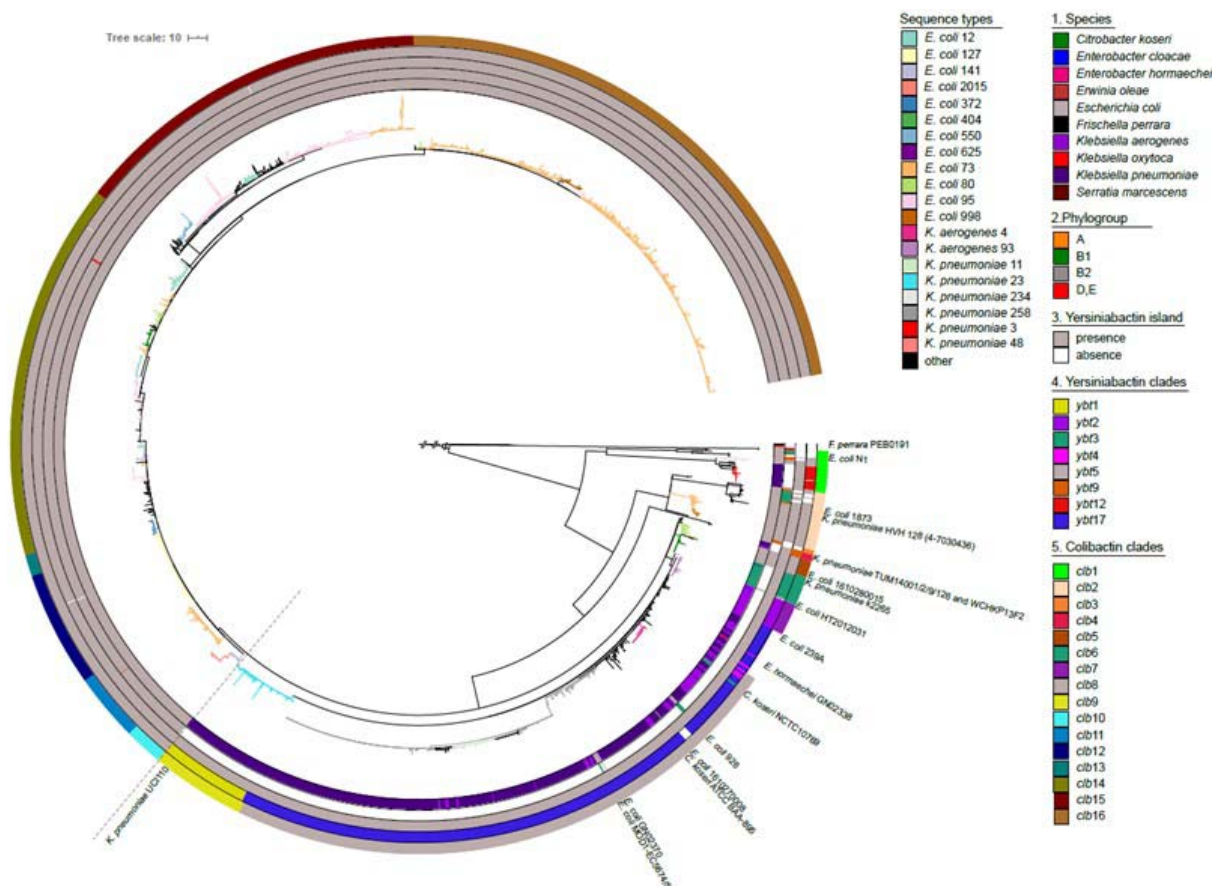


B)



748

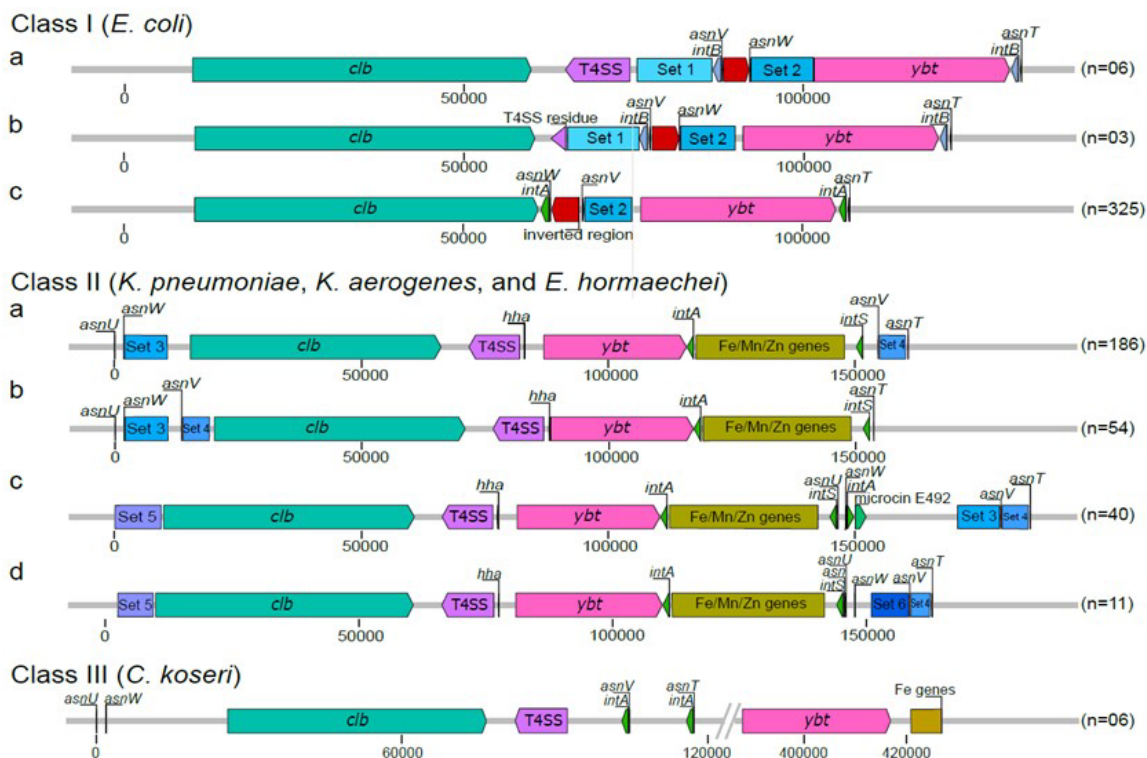
749 **Figure 2.** Maximum-likelihood based phylogenetic analysis of the colibactin and yersiniabactin  
 750 determinants. (A) Phylogenetic tree of the colibactin gene cluster (collapsed), B) phylogenetic  
 751 tree of the corresponding *ybt* determinants (collapsed) using the genetically distant *K.*  
 752 *michiganensis* strains as outgroup (Lam *et al.*, 2018). Additionally, the yersiniabactin sequence  
 753 type (YbST) as defined by Lam and colleagues (Lam *et al.*, 2018) associated with individual  
 754 bacterial clades are indicated. The branch colors in both trees depict the prominent bacterial  
 755 sequence type of the clade.



756

757 **Figure 3.** Maximum-likelihood based phylogeny of the colibactin gene cluster detected in 2169  
 758 enterobacterial genomes. Every leaf represents a single sequence variant of the *clb* gene  
 759 cluster, which can be allocated to different lineages and clades. From innermost to outermost,  
 760 the 1<sup>st</sup> circle indicates the species harboring the *clb* determinant; the 2<sup>nd</sup> circle shows the *E.*  
 761 *coli* phylogroup, the 3<sup>rd</sup> circle shows the presence/absence of the *ybt* operon; the 4<sup>th</sup> circle  
 762 shows the yersiniabactin sequence types (YbST) of the *ybt* determinant (from Fig. 1B) that  
 763 correspond to the *pks* island lineage present in the individual genome. The 5<sup>th</sup> circle shows the  
 764 different colibactin sequence types (CbST) of the *clb* gene cluster. The branch colors in the  
 765 center of the tree depict the prominent bacterial sequence types (Fig. 1). The large conserved  
 766 *E. coli* phylogroup B2 clade is separated from the large *Klebsiella* clade with a faint broken  
 767 line.

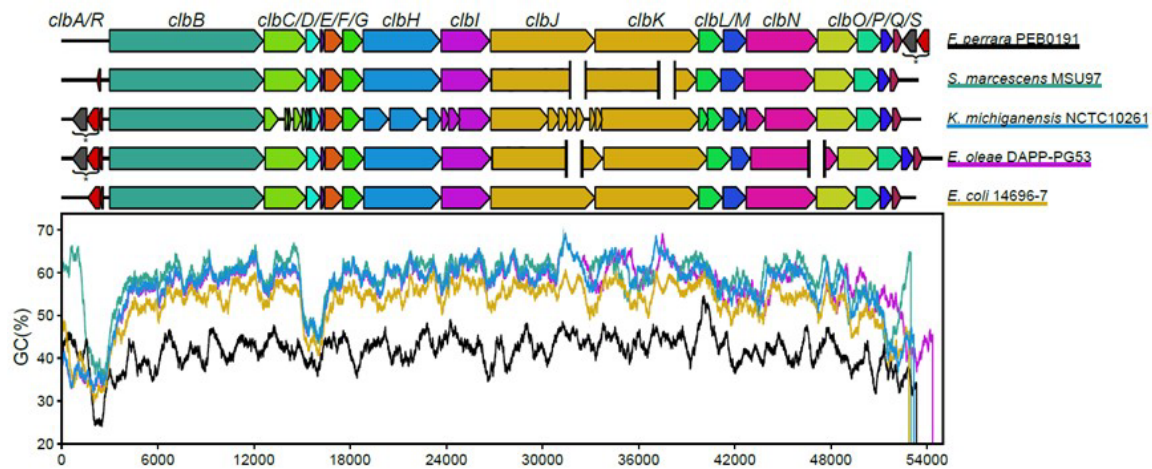
768



769

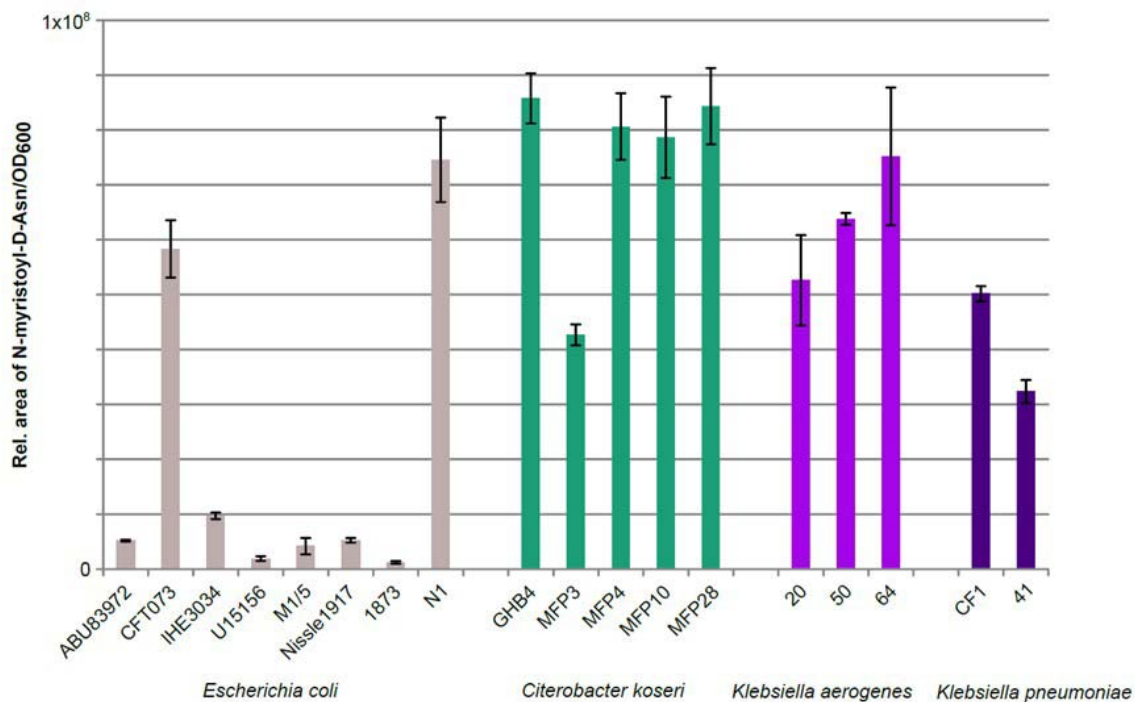
770 **Figure 4.** Structural variation of the colibactin and yersiniabactin-encoding chromosomal  
 771 region in *E. coli*, *E. hormaechei*, *K. pneumoniae*, *K. aerogenes* and *C. koseri*. The different  
 772 genetic structures and chromosomal insertion sites of the colibactin and/or yersiniabactin  
 773 determinants found within the three main structural classes are shown. The *clb* gene cluster  
 774 (teal green), T4SS module (purple), *ybt* gene cluster (pink), integrase genes (green), the  
 775 conserved sets of genes (Table S2) that are present up/downstream the two polyketide  
 776 determinants, classed into sets (blue boxes) and the Fe/Mn/Zn module (yellow) are shown.  
 777 The number of genomes included in the tested set of genomes that harbour the different  
 778 structural variants is indicated in brackets. The colibactin-yersiniabactin chromosomal regions  
 779 that do not conform to these major structures are as shown in Figure S4.

780



782 **Figure 5.** The structural organization and GC profiles of the *clb* determinants in the five most  
783 genetically distant bacterial strains (according to Fig. 2A). The genes that make up the  
784 homologous *pks* gene cluster found in *F. perrara* and the most distant *clb* determinants  
785 present in *S. marcescens*, *E. oleae*, *K. michiganensis* and *E. coli* (phylogroup E) strain 14696-7  
786 are depicted. The GC profile of the gene cluster in the different strains is shown alongside  
787 with the colors underlining the different species. SAM genes and *clbA* homologues (\*) are  
788 shown downstream of the *pks* gene cluster in *F. perrara* and upstream of the *clb* determinant  
789 in *K. michiganensis* and *E. oleae*. The gaps in assembly are shown with white spaces.

790

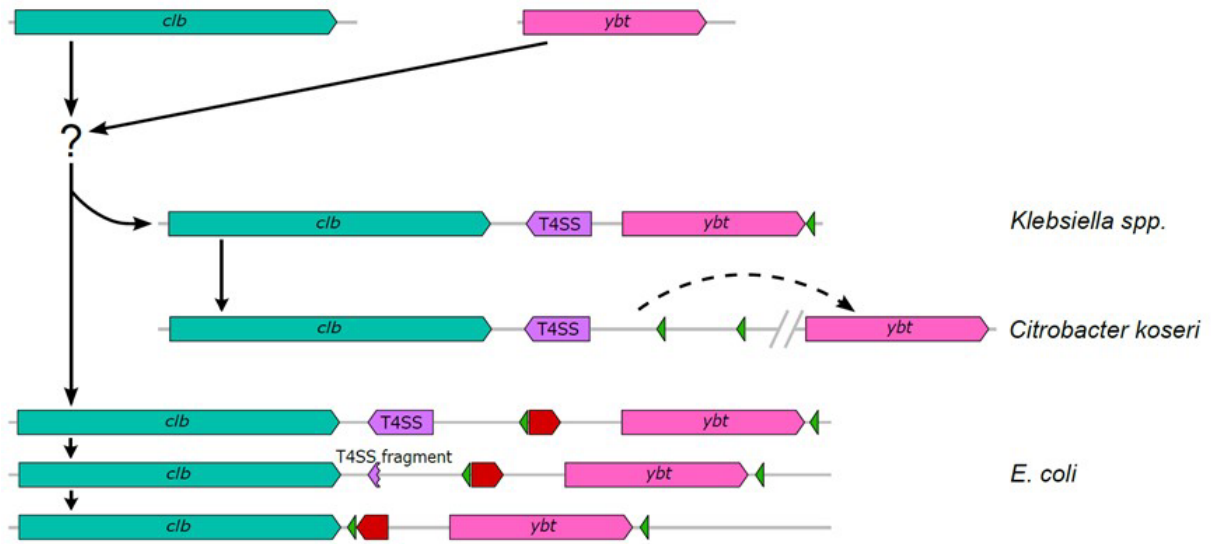


791

792 **Figure 6.** Comparison of colibactin production of different strains accessed by quantification  
793 of the precolibactin cleavage product N-Asn-D-myristol. This assay enabled us to compare the  
794 ability of different strains and different species to produce colibactin under controlled  
795 conditions *in vitro*. Measurements were conducted based on three biological replicates,  
796 means with standard deviations are as shown.

797





798

799 **Figure 7.** Schematic representation of the predicted evolution of the colibactin-yersiniabactin  
800 genomic region in *Enterobacteriaceae*. The different elements of this region, i.e. the *clb*  
801 determinant (teal green), T4SS module (purple), *ybt* gene cluster (pink), integrase genes  
802 (green), and an invertible subset of genes (red arrow) are shown. Based on available genome  
803 sequence data, we suggest a development from single MGEs containing the *clb* determinant  
804 and the *ybt* gene cluster, respectively, towards the structural arrangement of both polyketide  
805 determinants, which is now mainly found in enterobacterial populations. Black arrows (solid  
806 or dashed) indicate possible directions of development and DNA rearrangements. After the  
807 merge of the *clb* and *ybt* gene clusters into one MGE, represented by ICEKp10, there is  
808 evidence that three different structural variants have evolved from it: In *Klebsiella* spp.  
809 strains, the ICEKp10 has remained intact, whereas in *C. koseri* strains, a DNA rearrangement  
810 and re-localization of the *ybt* determinant to a different chromosomal position has taken  
811 place. In *E. coli*, a gradual loss of the T4SS module and the inversion of a gene set between the  
812 two polyketide determinants led to immobilisation or stabilisation of the ICE thus resulting the  
813 two pathogenicity islands known as *pks* island and HPI, respectively.

814

815

## 816 Supplemental Material

817 includes Supplementary Tables S1 (*clb*-positive bacterial strains), S2 (*clb*-positive STs in *E. coli*  
818 and *Klebsiella*), S3 (conserved gene sets associated with the *clb* and *ybt* determinants), S4 (list  
819 of prokaryotic species included into the *clb* screen), S5 (list of ClbSTs), and S6 (list of YbSTs)





**Table 1. Bacterial species tested positive for the presence of the colibactin determinant**

	No. of strains screened	<i>clb</i> -positive strains
<i>Escherichia coli</i>	19,183	1,462
<i>Klebsiella pneumoniae</i>	8,040	572
<i>Klebsiella aerogenes</i>	225	101
<i>Citrobacter koseri</i>	48	27
<i>Enterobacter hormaechei</i>	711	2
<i>Enterobacter cloacae</i>	610	2
<i>Serratia marcescens</i>	515	1
<i>Klebsiella michiganensis</i>	94	1
<i>Salmonella</i> sp.	8	1*
<i>Erwinia oleae</i>	1	1
<i>Frischella perrara</i>	3	3

\*unverified source organism (excluded from Refseq)