

1 **Supplemental information.**

2

3 **Gut symbionts from distinct hosts exhibit genotoxic**
4 **activity via divergent colibactin biosynthetic pathways**

5

6 Philipp Engel, Maria I. Vizcaino, and Jason M. Crawford

7

8 **Content**

9 Figures S1-S7

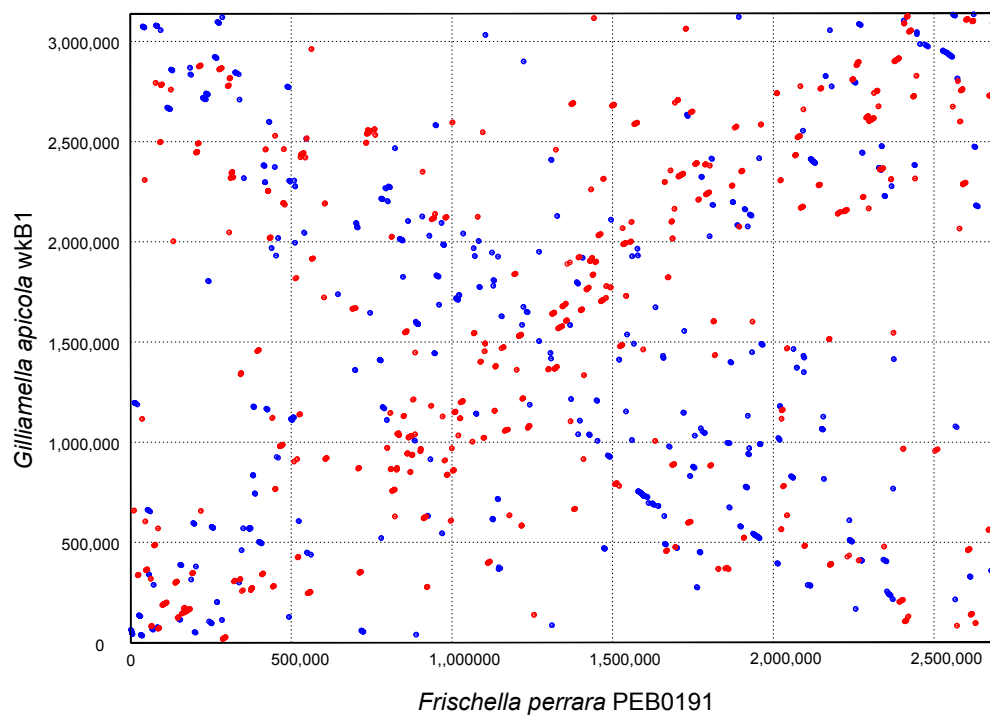
10 Tables S1-S5

11 Supplemental materials and methods

12 Supplemental references

13

14



15

16

17 **Figure S1.** Dotplot analysis of the genomes of *G. apicola* wkB1 and *F. perrara*

18 PEB0191. Plots were generated with Promer, a command of the software

19 package MUMmer v3.23 (1). Promer generates local alignments based on amino

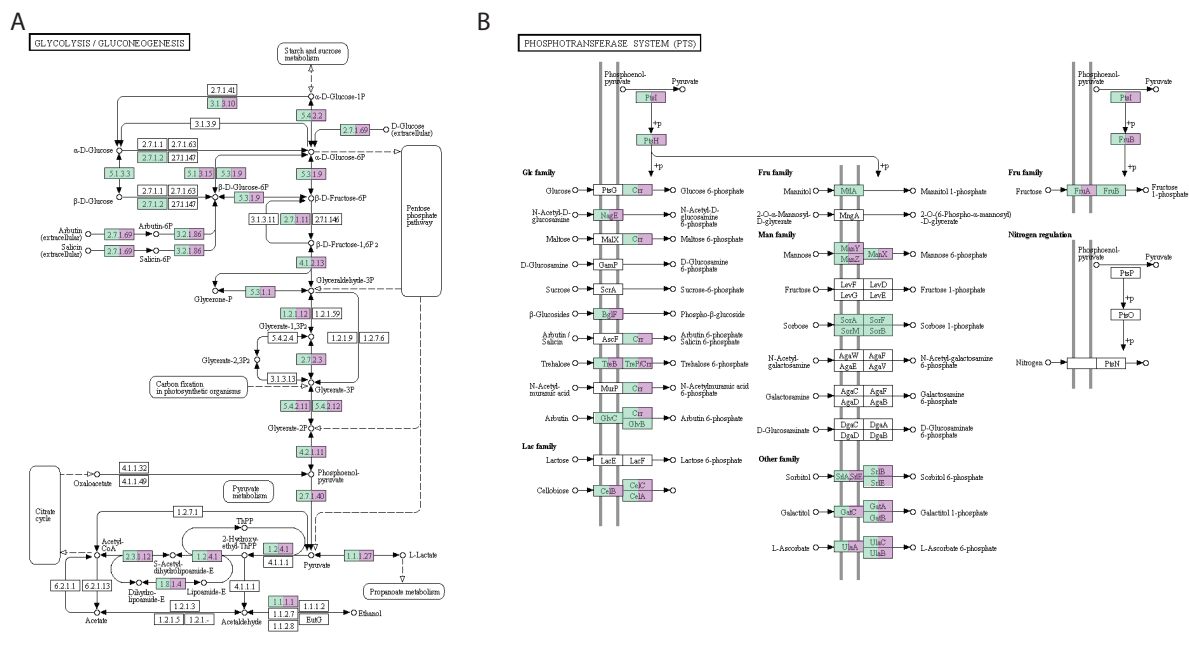
20 acid sequences and plots aligned regions onto the x-axis and y-axis representing

21 the genome positions of *F. perrara* PEB0191 and *G. apicola* wkB1, respectively.

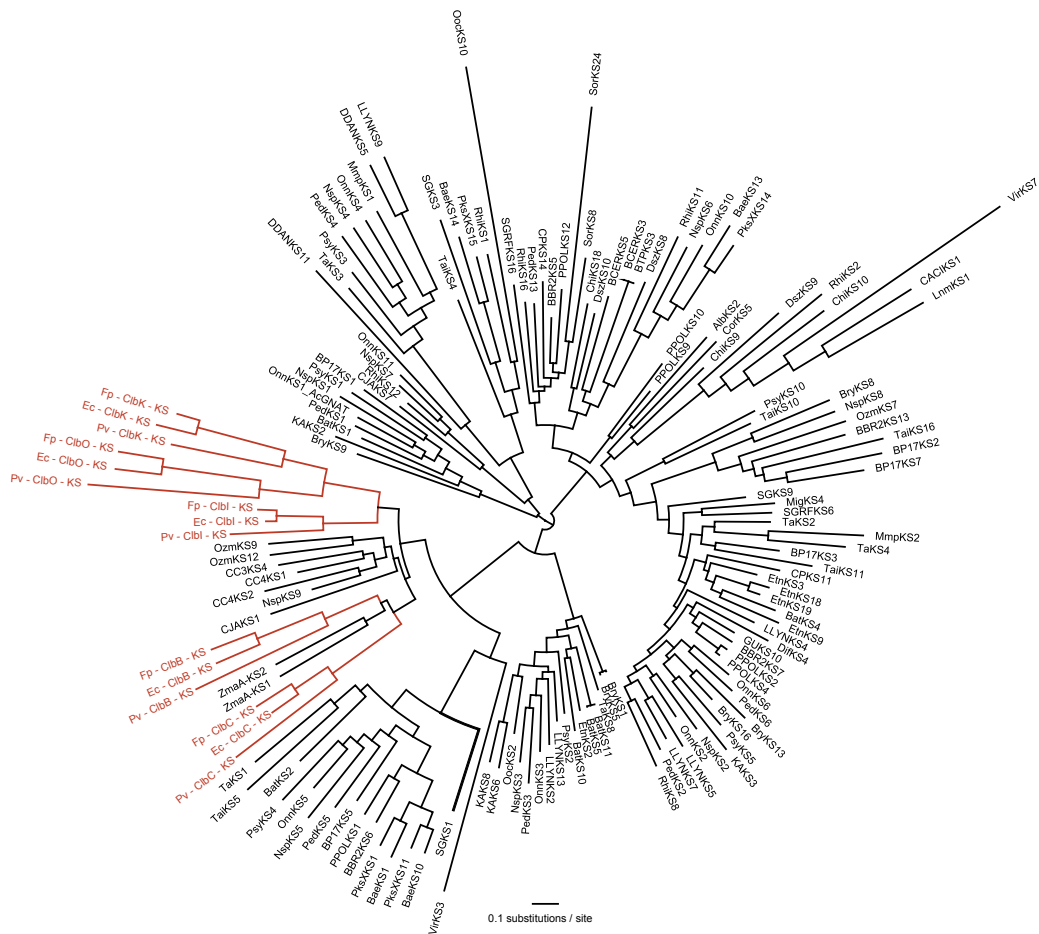
22 Default parameters were used. Red and blue dots represent alignments on the

23 same strand and on the opposite strand, respectively.

24



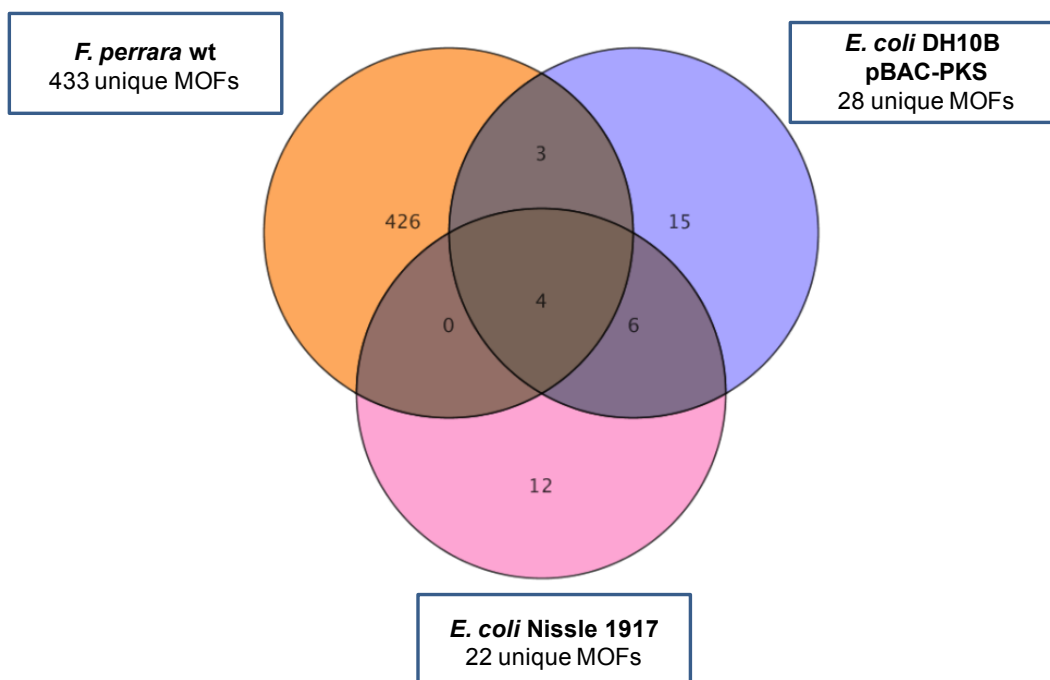
31
 32 **Figure S3.** Glycolysis (A) and phosphotransferase system (B) genes in the genomes of the two honey bee gut symbionts *F.*
 33 *perrara* and *G. apicola* wkB1. Gene functions identified in the genomes of *F. perrara* and *G. apicola* wkB1 are highlighted in
 34 magenta and green, respectively. Other gene functions are either absent or could not be identified.



36
 37 **Figure S4.** Maximum likelihood tree of ketosynthase (KS) domains of 154
 38 polyketide synthases (PKS) proteins. Amino acid sequences of KS domains were
 39 obtained from a previous study (2) and aligned together with the KS domains of
 40 ClbC, ClbK, and ClbO (*trans*-AT PKS), and ClbB and ClbI (*cis*-AT PKS) of *E. coli*
 41 IHE3034 (Ec), *F. perrara* (Fp), and *Pseudovibrio* FO-BEG1 (Pv). The phylogenetic
 42 tree was obtained with PhyML (3) using default parameters. KS domains of
 43 homologous Clb PKS proteins, highlighted in red, cluster together. The next most
 44 closely related KS domains are found in the following *cis*-AT PKS: ZmA of
 45 *Bacillus cereus* UW85 involved in the synthesis of zwittermicin (4), OzmK of
 46 *Streptomyces albus* JA3453 involved in the synthesis of oxazolomycin (5), NspD
 47 of *Nostoc* sp. 'Peltigera membranacea cyanobiont' involved in the synthesis of

48 nosperin (2), VirA of *Streptomyces virginiae* involved in the synthesis of
49 virginiamycin A (6), several PKSs of *Clostridium cellulolyticum* H10 (CC), and one
50 PKS of *Cellvibrio japonicus* Ueda107 (CJA).

51



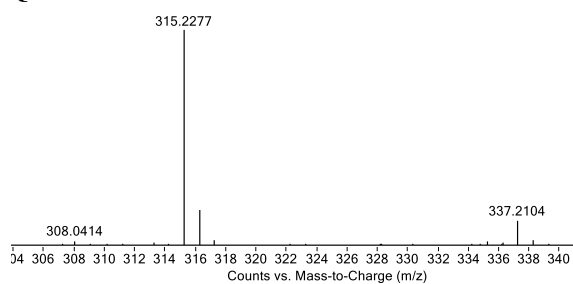
52

53

54 **Figure S5.** Venn Diagram comparison of the unique molecular features (MOFs)
 55 found in wild-type colibactin-bearing bacteria. The orange circle shows all MOFs
 56 of *F. perrara* PEB0191 (wt), from which MOFs found in the media controls were
 57 removed. The purple circle shows all MOFs of *E. coli* DH10B pBAC-PKS, from
 58 which MOFs found in *E. coli* DH10B pBAC-control were removed. The magenta
 59 circle shows all MOFs of *E. coli* Nissle 1917, from which MOFs found in *E. coli*
 60 Nissle 1917 Δclb (complete deletion of the genomic island) were removed. There
 61 are seven metabolites of *F. perrara*, shared among the different strains, four
 62 shared among all three (see Table S4).

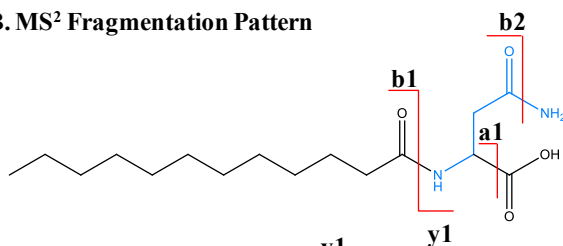
63

A. ESI-QTOF-HRMS

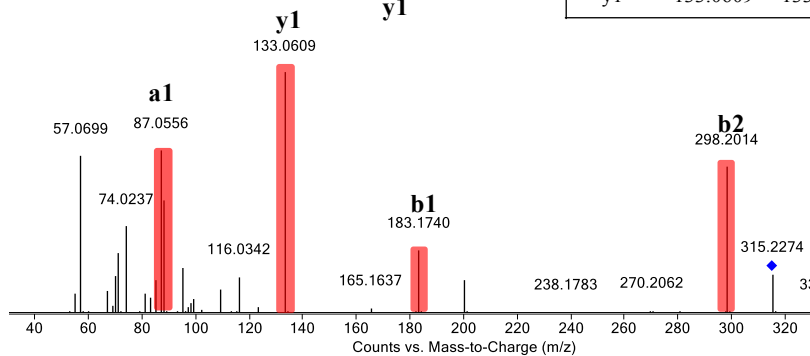


[M+H]⁺
Obs: 315.2277
Calc: 315.2284
Error (ppm): -2.22

B. MS² Fragmentation Pattern



Species	Obs. Mass	Calc. mass	Error (ppm)
a1	87.0556	87.0553	3.4
b1	183.1740	183.1743	1.6
b2	298.2014	298.2013	0.3
y1	133.0609	133.0608	0.7



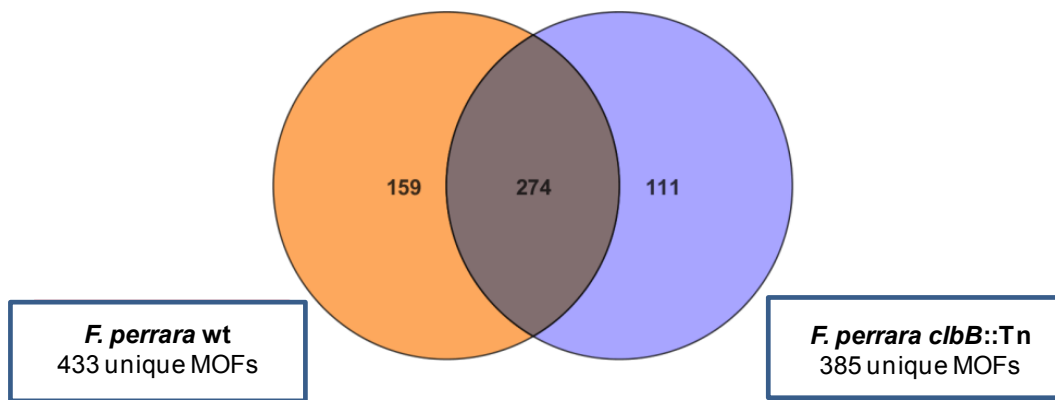
64

65

66 **Figure S6.** (A) ESI-QTOF-HRMS analysis and (B) MS² fragmentation pattern of *F.*
67 *perrara* metabolite **1**. Fragmented species shown in the table are highlighted in
68 the proposed structure of **1** and the corresponding MS² spectra. Observed mass
69 (Obs mass) and calculated mass (Calc. mass) are shown for both MS and MS²
70 ions, and all ppm errors were calculated to be less than 5.0.

71

72



73

74

75 **Figure S7.** Venn diagram comparison of the unique molecular features (MOFs)
76 found in *F. perrara wt* (PEB0191) and *F. perrara clbB::Tn*. MOFs found in the
77 media control were removed from both samples.

78

79 **Table S1.** *F. perrara*-specific genes are found in the xls document accompanying
80 this publication.
81

82 **Table S2.** Adenylation domain specificity predicted by three bioinformatic tools.

83

Adenylation domain^a	AntiSmash (Stachelhaus code)	PKS/NRPS Analysis Web-site	NRPSpredictor2^b
ClbB_Fp	valine	Pps3-M1-Glu	valine (70%)
ClbB_Pv	valine	no hit	valine (90%)
ClbB_Ec	valine	Pps3-M2-Val	valine (80%)
ClbH_Fp	serine	PhsB-M1-Ser	serine (100%)
	beta-hydroxy-tyrosine	no hit	valine (60%)
ClbG/H_Pv	serine	PhsB-M1-Ser	serine (100%)
	beta-hydroxy-tyrosine	no hit	valine (60%)
ClbH_Ec	serine	PhsB-M1-Ser	serine (100%)
	beta-hydroxy-tyrosine	no hit	valine (60%)
ClbJ_Fp	glycine	no hit	glycine (100%)
	cysteine	PchE-M1-Cys	cysteine (100%)
ClbJ_Pv	glycine	PchE-M1-Cys	glycine (100%)
	cysteine	NosC-M2-Gly	cysteine (100%)
ClbJ_Ec	glycine	no hit	glycine (100%)
	cysteine	MtaD-M1-Cys	cysteine (100%)
ClbK_Fp	cysteine	MtaD-M1-Cys	cysteine (100%)
ClbK_Pv	cysteine	PchE-M1-Cys	cysteine (100%)
ClbK_Ec	cysteine	MtaD-M1-Cys	cysteine (100%)
ClbN_Fp	asparagine	Cda2-M3-Asn	asparagine (100%)
ClbN_Pv	asparagine	Cda2-M3-Asn	asparagine (100%)
ClbN_Ec	asparagine	Cda2-M3-Asn	asparagine (100%)

84 ^aFp, *F. perrara*, Pv, *Pseudovibrio* sp., Ec, *E. coli* IHE3934

85 ^bPercentages indicate confidence scores

86

87 **Table S3.** Secondary structure analysis with Phyre2 (8) identifies inactivated
 88 acyltransferase (AT) domains encoded in PKS genes *clbC*, *clbK*, and *clbO*. *ClbG* is a
 89 conserved trans-AT PKS, which served as a positive control. For each query, only
 90 the three best hits are shown.

<i>F. perrara</i> PEB0191	Rank	Hit	Confidence %	^a Coverage aa	^b Gaps aa	^c Motifs (query / template)
ClbC 435 aa - 786 aa	1	Transferase of PKS: c3tzzA	100%	298	149	GEGLG / GQ SLG
	2	Transferase of PKS: c2qo3A	100%	308	140	GEGLG / GH SQ G
	3	Transferase of PKS: c2hg4A	100%	304	140	GEGLG / GH SQ G
ClbK 421 aa - 724 aa	1	Hydrolase of PKS: c4oqjA	100%	276	6	GQGDG / ADRTE
	2	Transferase of PKS: c4mz0B	100%	274	125	GQGDG / GH SV G
	3	Transferase of PKS: c2hg4A	100%	275	184	GQGDG / GH SQ G
ClbO 427 aa - 748 aa	1	Transferase of PKS: c3tzzA	100%	289	155	GNQAG / GQ SL G
	2	Transferase of PKS: c2qo3A	100%	296	149	GNQAG / GH SQ G
	3	Transferase of PKS: c2hg4A	100%	297	144	GNQAG / GH SQ G
ClbG 1 aa - 413 aa	1	Transferase of PKS: c2qo3A	100%	354	10	GH SL G / GH SQ G
	2	Transferase of malonyl-CoA-ACP transacylase: c3eenA	100%	315	0	GH SL G / GH SL G
	3	Transferase of malonyl-CoA-ACP transacylase: c2qj3B	100%	314	0	GH SL G / GH SV G
<i>E. coli</i> IHE3034	Rank	Hit	Confidence %	Coverage aa	Gaps aa	Motifs (query / template)
ClbC 435 aa - 787 aa	1	Transferase of PKS: c4mz0B	98.70%	298	78	GAGTG / GH SV G
	2	Hydrolase of PKS: c4oqjA	98.20%	306	12	GAGTG / ADRTE
	3	Transferase of PKS: c4na3A	97.40%	301	5	GAGTG / -
ClbK 420 aa - 720 aa	1	Hydrolase of PKS: c4oqjA	100%	266	14	GDGDG / ADRTE
	2	Transferase of PKS: c2qo3A	100%	265	187	GDGDG / GH SQ G
	3	Transferase of fatty acid synthetase: c2vz8A	100%	286	170	GDGDG / GH SL G
ClbO 427 aa - 746 aa	1	Transferase of fatty acid synthetase: c2vz8A	100%	306	147	GYLTG / GH SL G
	2	Transferase of PKS: c3tzzA	100%	316	58	GYLTG / GQ SL G
	3	Transferase of PKS: c2qo3A	100%	296	150	GYLTG / GH SQ G
ClbG 1 aa - 422 aa	1	Transferase of PKS: c2qo3A	100%	349	12	GH SL G / GH SQ G
	2	Transferase of PKS: c2hg4A	100%	351	13	GH SL G / GH SQ G
	3	Transferase of PKS: c3tzzA	100%	367	15	GH SL G / GQ SL G
<i>Pseudovibrio</i> FO-BEG1	Rank	Hit	Confidence	Coverage aa	Gaps aa	Motifs(query / template)
ClbC 434 aa - 752 aa	1	Transferase of PKS: c3tzzA	100%	271	183	GQGDG / GQ SL G
	2	Transferase of fatty acid synthetase: c2vz8A	100%	293	167	GQGDG / GH SL G
	3	Transferase of PKS: c2qo3A	100%	280	187	GQGDG / GH SQ G
ClbK 426 aa - 744 aa	1	Hydrolase of PKS: c4oqjA	100%	279	11	FAGNK / ADRTE
	2	Transferase of PKS: c3tzzA	100%	278	172	FAGNK / GQ SL G
	3	Transferase of PKS: c2qo3A	100%	278	166	FAGNK / GH SQ G
ClbO 444 aa - 799 aa	1	Transferase of fatty acid synthetase: c2jfkD	97.9%	176	14	GDGTG / GH SL G
	2	Transferase of PKS: c4mz0B	97.3%	287	59	GDGTG / GH SV G
	3	Transferase of PKS: c3tzzA	97.2%	188	28	GDGTG / GQ SL G
ClbG 1 aa - 201 aa	1	Transferase of malonyl-CoA-ACP transacylase: c3tqeA	100%	200	0	GH SL G / GH SL G
	2	Transferase of malonyl-CoA-ACP transacylase: c3eenA	100%	200	0	GH SL G / GH SL G
	3	Transferase of malonyl-CoA-ACP transacylase: c3ptwA	100%	200	0	GH SL G / GL SL G

91 ^aamino acids (aa) of the query sequence in the alignment

92 ^bnumber of alignment gaps in the query sequence given in amino acids (aa)

93 ^caligned sequence motifs are shown with the conserved active site Serine residue highlighted in red

94 **Table S4.** Seven colibactin-pathway dependent metabolites shared among *F.*
 95 *perrara* and two colibactin-bearing *E. coli* strains.^a

96

[M]	[M+H] ⁺	RT	Fp wt	Fp <i>clbB</i> ::Tn	pBAC-PKS	EcN wt	Proposed structure
207.0568	208.0640	1.18	8.E+05	(4/5)	4.E+05	(1/5)	
258.1221	259.1293	15.35	8.E+05	-	3.E+05	(3/5)	
314.2209	315.2281	14.44	2.E+08	1.E+06	2.E+06	8.E+05	1
340.2366	341.2440	15.34	7.E+06	-	3.E+06	1.E+06	3
342.2520	343.2593	16.67	2.E+07	(2/5)	1.E+08	4.E+07	2
368.2678	369.2749	17.37	5.E+05	-	5.E+05	4.E+05	4
370.2839	371.2903	18.95	9.E+05	-	3.E+05	(3/5)	5

97

98 ^aMasses not detected in all five biological replicates were represented by the
 99 number of times they were observed in five replicates). The symbol, -, indicates
 100 that the mass was not observed in any biological replicates. Abbreviations are as
 101 follows: Fp wt, *F. perrara* PEB0191, Fp *clbB*::Tn, *F. perrara clbB*::Tn, pBAC-PKS, *E.*
 102 *coli* DH10B pBAC-PKS, EcN wt, *E. coli* Nissle 1917.

103

104 **Table S5.** Colibactin-pathway dependent metabolites from *F. perrara*.

[M]	<i>m/z</i>	<i>m/z</i> ^a	RT (min)	Fp <i>wt</i>	Fp <i>clbB</i> ::Tn	Proposed structure ^b
227.0962 ^{c,d}	228.1040	n.d.	14.450	3E+5	0	
258.1221	259.1293	1	15.350	8E+5	0	
286.1897	287.1970	1	12.128	1E+6	0	6
297.1947	298.2020	1	14.443	3E+6	0	
312.2058	313.2129	1	13.382	1E+6	0	7
312.2067	313.2140	1	13.655	3E+5	0	
313.0199	314.0270	1	5.648	9E+5	0	
314.2209	315.2281	1	14.444	2E+8	1E+6	1
328.2364	329.2443	1	15.557	4E+6	0	8
340.2366	341.2440	1	15.345	7E+6	0	3
341.1810 ^{b,c}	683.3676	2	14.442	3E+5	0	
342.2520	343.2593	1	16.673	2E+7	(2/5)	2
342.9579 ^c	343.9651	n.d	11.335	6E+5	0	
368.2678	369.2749	1	17.372	5E+5	0	4
370.2839	371.2903	1	18.953	9E+5	0	5
414.1452 ^c	415.1529	nd	14.453	3E+5	0	
434.2220	218.1183	2	14.447	2E+6	0	
666.3902	334.2025	2	14.448	2E+6	0	
980.6098	491.3124; 981.6147	2; 1	14.443	2E+6	0	
996.5832 ^c	997.5929	1	14.444	5E+5	0	

105

106 ^an.d. = could not be determined based on low abundance

107 ^bSee Fig. 4

108 ^cRaw abundance determined from diluted metabolomics samples

109 ^dMS² fragmentation was not successfully acquired

110

111

112

113 **Supplemental materials and methods**

114 **Genome sequencing, assembly, and annotation.** Genomic DNA was isolated
115 from *F. perrara* PEB0191 grown for 2 days after re-streaking using the phenol-
116 chloroform method. Libraries for SMRT sequencing (Pacific Biosciences) were
117 constructed as recommended by the manufacturer. Sequencing of two SMRT
118 cells was necessary to obtain 64,460 quality-filtered reads of 2.9 kb average
119 length. Error correction, assembly, and consensus sequence polishing were
120 carried out with the PacBio HGAP pipeline (9) resulting in two contigs, one large
121 contig that covered the entire *F. perrara* genome and a second, smaller contig,
122 which was removed from the assembly due to its low read coverage. The
123 Illumina paired-end library with approximate insert sizes of 400 bp was
124 constructed from the genomic DNA following Illumina standard protocols for
125 genome sequencing using four PCR amplification cycles with the Bio HiFi
126 polymerase (Kapa Biosystems, Woburn, MA, USA). Illumina sequencing was
127 carried out on a HiSeq2000 machine in a single 2 × 100 bp lane at the Yale Center
128 for Genome Analysis. Illumina reads were trimmed on quality with CLC
129 Genomics Workbench (CLC Bio) as previously published (10). In total, 5,411,774
130 reads passed the quality filter. These Illumina reads were mapped with BWA
131 (11) against the PacBio assembly and inspected for misassemblies using the
132 sequence assembly viewer Tablet (12). A small number of sequencing errors and
133 one misassembly were detected and corrected. The modified assembly was again
134 verified by read mapping and by visual inspection. Overhanging ends were
135 trimmed and the chromosome position 1 set to the origin of replication, which
136 was detected using the online tool Ori-Finder (13).

137

138 **Transposon mutagenesis.** A transposon mutant library of about ~500 clones
139 was generated by conjugation of *F. perrara* PEB0191 with *E. coli* β 2163
140 harboring plasmid pBT20. Cells of *F. perrara* PEB0191 grown for 30 h were
141 harvested from two GMM (gut microbiota medium) (14) agar plates and
142 resuspended in 100 μ L of 1x PBS. An overnight culture of *E. coli* β 2163 harboring
143 plasmid pBT20 was grown for 8 h in LB supplemented with 30 μ g/mL
144 gentamicin, 100 μ g/mL ampicillin, and 0.3 mM diaminopimelic acid (DAP),
145 washed once in 1x PBS, and resuspended in 1 mL 1x PBS. 20 μ L of resuspended
146 *E. coli* cells were mixed with the resuspended *F. perrara* cells, spread onto a
147 nitrocellulose filter (Sartorius) on a GMM agar plate supplemented with 0.1 mM
148 DAP, and incubated for 16 h at 37° C in a 5% CO₂ incubator. Bacteria were
149 recovered in 1 mL 1x PBS, washed once, and distributed on eight GMM agar
150 plates supplemented with 12.5 μ g/mL gentamicin. After three days of incubation
151 under anaerobic conditions at 37° C, single colonies (transposon mutants) were
152 picked, resuspended in 150 μ L GMM supplemented with 12.5 μ g/mL gentamicin
153 and grown at 37°C under anaerobic conditions in a 96-well plate. Pools of 96 and
154 48 transposon mutants were generated, and genomic DNA was isolated using the
155 DNAeasy kit (Qiagen). The remaining bacterial suspensions of the single
156 transposon mutants were frozen at -80° C after adding 30 μ L of glycerol. We
157 screened all mutant pools by PCR with nine different primer combinations, each
158 combination consisting of two outward-facing primers (prRND1 and
159 prRND1rev) annealing to the two ends of the transposon and one of nine
160 primers (prPE209-217) annealing to different regions of the *clb* GI. PCRs were
161 analyzed by agarose gel electrophoresis to identify amplicons possibly
162 originating from a transposon insertion into the *clb* GI. Single transposon
163 mutants of pools suspected to contain an insertion mutant in the *clb* GI were

164 then individually screened by PCR. The transposon integration in the gene *clbB*
165 was confirmed by PCR over the integration site (prPE245 and prPE246) and by
166 Sanger sequencing. All strains, plasmids, and primers used in this study are
167 summarized in Table 1.

168

169 **Analysis of γ -H2AX phosphorylation levels in HeLa cells.** γ -H2AX
170 phosphorylation levels in HeLa cells were analyzed to detect the activation of a
171 DNA damage response. Therefore, 3×10^5 cells were seeded in each well of a 6-
172 well plate and incubated for 6 h allowing cells to adhere. Then, cells were
173 transiently infected with bacterial cultures for 4 h at the indicated multiplicity of
174 infection (MOI). Bacteria were removed by washing the cells 3-6 times with
175 DMEM/5% FCS. Subsequently, HeLa cells were incubated for 12h in DMEM/5%
176 FCS supplemented with 200 μ g/mL gentamicin. Then, cells were removed from
177 the culture dish with 100 μ L 0.025% trypsin/0.01% EDTA and spun down in 1
178 mL cell medium in a microcentrifuge at 200 g. After washing in 1x PBS, cells
179 were fixed in 4% paraformaldehyde for 10 min at room temperature (RT), then
180 incubated in 20 mM NH_4Cl for 2 min, and washed again in 1x PBS. HeLa cells
181 were resuspended in 100 μ L 1x PBS, and 900 μ L methanol was added during
182 constant vortexing. After incubation for 30 min on ice, cells were washed in 0.5%
183 BSA and incubated in 200 μ L 0.5 % BSA supplemented with 0.5 μ L anti- γ -H2AX
184 primary antibody (clone 20E3, Cell Signaling) for 1h at RT. Cells were again
185 washed in 0.5 % BSA and then incubated in 100 μ L 0.5 % BSA supplemented
186 with 2 μ L of the secondary antibody conjugated to FITC (goat anti-rabbit
187 AB97199, ABCAM) for 30 min in the dark at RT. Cells were washed once in 0.5%
188 BSA, resuspended in 500 μ L 1x PBS, and then analyzed by flow cytometry using a
189 FACSVerseTm flow cytometer from BD Bioscience.

190 **Supplemental references**

- 191 **1. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al.**
192 2004. Versatile and open software for comparing large genomes. *Genome Biol*
193 **5**:R12.
- 194 **2. Kampa A, Gagunashvili AN, Gulder TAM, Morinaka BI, Daolio C,**
195 **Godejohann M, et al.** 2013. Metagenomic natural product discovery in lichen
196 provides evidence for a family of biosynthetic pathways in diverse symbioses.
197 *Proc Natl Acad Sci USA* **110**:E3129–37.
- 198 **3. Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O.**
199 2010. New algorithms and methods to estimate maximum-likelihood
200 phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* **59**:307–321.
- 201 **4. Kevany BM, Rasko DA, Thomas MG.** 2009. Characterization of the complete
202 zwittermicin A biosynthesis gene cluster from *Bacillus cereus*. *Appl Environ*
203 *Microbiol* **75**:1144–1155.
- 204 **5. Zhao C, Coughlin JM, Ju J, Zhu D, Wendt-Pienkowski E, Zhou X, et al.** 2010.
205 Oxazolomycin biosynthesis in *Streptomyces albus* JA3453 featuring an
206 ‘acyltransferase-less’ type I polyketide synthase that incorporates two distinct
207 extender units. *J Biol Chem* **285**:20097–20108.
- 208 **6. Pulsawat N, Kitani S, Nihira T.** 2007. Characterization of biosynthetic gene
209 cluster for the production of virginiamycin M, a streptogramin type A antibiotic,
210 in *Streptomyces virginiae*. *Gene* **393**:31–42.
- 211 **7. Engel P, Kwong WK, Moran NA.** 2013. *Frischella perrara* gen. nov., sp.
212 **nov., a gammaproteobacterium isolated from the gut of the honeybee, *Apis***

- 213 *mellifera*. Int J Syst Evol Microbiol **63**:3646–3651.
- 214 8. **Kelley LA, Sternberg MJE**. 2009. Protein structure prediction on the Web: a
215 case study using the Phyre server. Nat Protoc **4**:363–371.
- 216 9. **Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al**.
217 2013. Nonhybrid, finished microbial genome assemblies from long-read SMRT
218 sequencing data. Nat Methods **10**: 563-569.
- 219 10. **Engel P, Martinson VG, Moran NA**. 2012. Functional diversity within the
220 simple gut microbiota of the honey bee. Proc Natl Acad Sci USA **109**:11002–
221 11007.
- 222 11. **Li H, Durbin R**. 2009. Fast and accurate short read alignment with Burrows-
223 Wheeler Transform. Bioinformatics **25**:1754-60.
- 224 12. **Milne I, Bayer M, Cardle L, Shaw P, Stephen G, Wright F, Marshall D**.
225 2010. Tablet--next generation sequence assembly visualization. Bioinformatics
226 **26**:401-402.
- 227 13. **Gao F, Zhang C-T**. 2008. Ori-Finder: A web-based system for finding oriCs in
228 unannotated bacterial genomes. BMC Bioinformatics **9**:79.
- 229 14. **Goodman AL, Kallstrom G, Faith JJ, Reyes A, Moore A, Dantas G, et al**.
230 2011. Extensive personal human gut microbiota culture collections characterized
231 and manipulated in gnotobiotic mice. Proc Natl Acad Sci USA **108**:6252–6257.