|   |  | A101555A_nr-  |
|---|--|---------------|
|   |  | reporting-    |
|   |  | summary_2020_ |
|   |  | 04_22         |
|   |  | updated.pdf   |
| 4 |  |               |

A. Additional Supplementary Files

| Туре                | Number<br>If there are multiple files<br>of the same type this<br>should be the numerical<br>indicator. i.e. "1" for<br>Video 1, "2" for Video 2,<br>etc. | Filename<br>This should be the name<br>the file is saved as when<br>it is uploaded to our<br>system, and should<br>include the file<br>extension. i.e.: Smith_<br>Supplementary_Video_1.<br>mov | Legend or<br>Descriptive Caption<br>Describe the contents of<br>the file |
|---------------------|---|---|--|
| Supplementary Table | 1   | Tables.xlsx   | Supplementary<br>Tables 1-2  |

| 8  |   |
|----|---|
| 9  | Colibactin DNA damage signature indicates mutational impact in colorectal cancer  |
| 10 |   |
| 11 | Paulina J. Dziubańska-Kusibab <sup>1†</sup> , Hilmar Berger <sup>1†</sup> , Federica Battistini <sup>2#</sup> , Britta A.M.                               |
| 12 | Bouwman <sup>3#</sup> , Amina Iftekhar <sup>1#</sup> , Riku Katainen <sup>4#</sup> , Tatiana Cajuso <sup>4</sup> , Nicola Crosetto <sup>3</sup> , Modesto |
| 13 | Orozco <sup>2,5</sup> , Lauri A. Aaltonen <sup>4</sup> , and Thomas F. Meyer <sup>1*</sup>  |
| 14 |   |
| 15 | + Shared first authors  |
| 16 | <sup>#</sup> Shared second authors  |
| 17 | <sup>1</sup> Department of Molecular Biology, Max Planck Institute for Infection Biology, 10117 Berlin,   |
| 18 | Germany   |
| 19 | <sup>2</sup> Institute for Research in Biomedicine (IRB Barcelona). The Barcelona Institute of Science  |
| 20 | and Technology, 08028 Barcelona, Spain  |
| 21 | <sup>3</sup> Science for Life Laboratory, Department of Medical Biochemistry and Biophysics, Karolinska   |
| 22 | Institutet, 17165 Stockholm, Sweden   |

- <sup>4</sup> Applied Tumor Genomics Research Program and Department of Medical and Clinical
- 24 Genetics, Medicum, University of Helsinki, 00014 Helsinki, Finland
- <sup>5</sup> Department of Biochemistry and Biomedicine, University of Barcelona, Barcelona, Spain
- 26
- 27 \* Corresponding author:
- 28 Prof. Thomas F. Meyer
- 29 Department of Molecular Biology
- 30 Max Planck Institute for Infection Biology
- 31 Charitéplatz 1
- 32 10117 Berlin
- 33 Germany
- 34 Email: <u>tfm@mpiib-berlin.mpg.de</u>
- 35

Abstract. The mucosal epithelium is a common target of damage by chronic bacterial 36 infections and accompanying toxins and most cancers originate from this tissue. We 37 investigated if colibactin, a potent genotoxin<sup>1</sup> associated with certain strains of 38 Escherichia coli<sup>2</sup>, creates a specific DNA damage signature in infected human colorectal 39 cells. Notably, the genomic contexts of colibactin-induced DNA double-strand breaks 40 (DSBs) were enriched for an AT-rich hexameric sequence motif, associated with distinct 41 42 DNA shape characteristics. A survey of somatic mutations at colibactin target sites of several thousand cancer genomes revealed significant enrichment of this motif in 43 colorectal cancers (CRCs). Moreover, the exact DSB loci corresponded with mutational hot 44 spots in cancer genomes, reminiscent of a trinucleotide signature previously identified in 45 healthy colorectal epithelial cells<sup>3</sup>. This work provides evidence for the etiological role of 46 colibactin in human cancer. 47

In humans, besides the prototypical role of *Helicobacter pylori* in gastric cancer<sup>4</sup>, several 48 bacterial species have been proposed to have a tumorigenic role in CRC, including 49 Fusobacterium nucleatum <sup>5</sup> and colibactin-producing strains of *E. coli* <sup>6,7</sup>. Mechanistic studies 50 51 have indicated distinct bacteria-induced cancer-promoting features, including the activation of inflammatory and growth-promoting signaling hubs and the induction of DNA damage<sup>8</sup>. 52 53 Genotoxins are widespread amongst bacterial species, where they participate in intermicrobial competition<sup>9</sup>. Several bacterial genotoxins, such as DNA alkylators duocarmycin<sup>10</sup>, 54 vatakemycin<sup>11</sup> and CC-1065 as well as polymerase inhibitors distamycin<sup>12</sup> and netropsin<sup>13</sup>, 55 show preferential binding to AT-rich motifs in the minor groove <sup>14</sup>. The B2 phylogenetic 56 57 group of *E. coli* harbors a 54 kilobase (kb)-long *pks* genomic island that encodes a polyketidepeptide hybrid, responsible for the production of the secondary metabolite colibactin<sup>1</sup>. This 58 59 toxin has been shown to induce double-strand breaks (DSBs) in host cells, followed by the activation of the G2-M DNA damage checkpoint pathway<sup>1</sup>. The chemical structure of several 60 colibactin precursors <sup>15-18</sup> and the isolation of colibactin-dependent N3 adenine adducts from 61 host DNA <sup>19</sup> have been recently described. The high-resolution mature structure of colibactin 62 63 depicts a highly symmetrical molecule that contains identical cyclopropane warheads at each end, eliciting DNA cross-links<sup>20</sup>. These DNA cross-links may be processed by depurination at 64 sites of colibactin action <sup>21</sup>. 65

66 We set out to explore the potential specificity of colibactin-induced DNA damage by 67 determining the characteristics of DNA breakpoints in infected human colorectal cells.

Interestingly, using a multimodal approach (Extended Data Fig. 1), we identified an infection
signature that was preserved in the genomes of about 10% of human CRCs.

70 To induce damage by colibactin, we infected human colorectal adenocarcinoma Caco-2 cells 71 with *pks+ or pks- E. coli* for 3 hrs, followed by fixation of the cells (Extended Data Fig. 2A). Controls included no treatment or treatment with DSB-inducing, topoisomerase II inhibitor 72 73 etoposide (Fig. 1A). Following infection, colibactin-induced DSBs were identified using sBLISS<sup>22</sup>, a sequencing-based DSB capture method, enabling unbiased identification of DSBs 74 75 at a genome-wide scale with nucleotide resolution (Fig. 1A and Methods). To verify that sBLISS was able to capture known patterns of etoposide-induced damage, the breakpoint 76 density around transcription start sites (TSS) was determined <sup>22-24</sup>. Indeed, an expected 77 78 increase in breakpoint counts around TSSs was observed in the etoposide control (Extended 79 Data Fig. 2B). Interestingly, unlike etoposide-induced DSBs, or those detected in the negative controls, DSBs induced upon *pks*+ *E. coli* infection failed to show any strong correlation with 80 known particular genomic regions (Extended Data Fig. 2C) using Locus Overlap Analysis 81 (LOLA)<sup>25</sup>. 82

83 While we hypothesized that sBLISS captured the preferred target sites of colibactin, we analyzed the nucleotide sequence context of induced DSBs (+/-10 nt) in pks+ E.coli- vs. pks-84 *E.coli*-infected cells, using discriminative motif discovery (DREME) <sup>26</sup> (Fig. 1B). One of these 85 86 motifs, a heptamer sequence, was strongly enriched in all replicates, overlapped the break point position, and included the consensus pattern AAWWTT with a 3' variable, single 87 nucleotide extension (Extended Data Fig. 3A). No enrichment for AAWWTT, referred to as 88 89 the colibactin damage motif (CDM) could be identified when comparing etoposide-induced 90 DSBs to untreated controls (Extended Data Fig. 3B). In order to evaluate how strong was the enrichment of this motif over all other motifs, we compared the enrichment of short 91 92 oligonucleotide sequences (i.e., pentanucleotides or hexanucleotides) around break sites following different treatments (Methods; Fig. 1C). We found that DSBs of pks+ E. coli-93 infected cells were generally enriched in AT-rich regions. This enrichment was particularly 94 high for the hexanucleotides AAATTT, AAAATT and AATATT, together with their 95 96 complementary mates. These sequence preferences of DSBs in pks+ E. coli-infected cells 97 were evident in comparison to both pks- E. coli-infected (Fig. 1D, top panel) and untreated 98 cells (Extended Data Fig. 4A). This enrichment was independently detected with almost identical relative values in all four biological replicates (Extended Data Table S1). 99

Importantly, no meaningful sequence enrichments were observed in the comparison 100 101 between pks- E. coli-infected cells and untreated cells (Fig. 1D, bottom panel). Hence, the 102 preference for AT-rich sequences is directly linked to the action of colibactin, rather than to 103 non-toxigenic E. coli infection. Similar enrichments were found for pentanucleotide motifs, 104 whereby sub-patterns of AAWWTT were most strongly enriched (Extended Data Fig. 4B). 105 Poly(dA:dT) tracts in eukaryotic genomes have been experimentally associated with genomic fragility <sup>27</sup>. Importantly, we found that the enrichment of other AT-rich sequences was 106 substantially weaker than that of AAWWTT (Extended Data Fig. 4C). Moreover, AT-rich 107 nucleosome-free regions<sup>28,29</sup> failed to show a distinct preference for AAWWTT, indicating 108 109 that the motif preference is not driven by increased access of colibactin to unprotected DNA regions (Extended Data Fig. 4D). 110

111 Binding of small molecules to DNA is partially determined by DNA shape characteristics (reviewed by Tse et al. <sup>30</sup>). We predicted the DNA shape parameters for the central positions 112 113 of all possible 4096 hexanucleotides and correlated them with the log2-ratios of hexanucleotide sequence enrichment at DSB positions in pks+ E. coli- vs. pks- E. coli-infected 114 115 cells. Remarkably, two out of the three colibactin-specific hexanucleotide sequences 116 (AAAATT/AATTTT and AAATTT) exhibited the narrowest minor groove widths as well as the 117 most extreme negative electrostatic potential (Fig. 2A, Extended Data Fig. 5A). Values for DNA stiffness descriptor <sup>31</sup> (k tot, see Methods) revealed that these tracts also possessed 118 119 high intrinsic rigidity (Supplementary Table S1), rendering them inert to distortion. Multivariate analysis of DNA shape parameters along the full length of all possible 120 121 hexanucleotides confirmed these extreme characteristics of CDM-related sequences (Fig. 2B). 122

123 In order to elucidate the binding between the DNA and colibactin, we built a molecular model using quantum mechanics calculations as a first structural prediction (see Methods). 124 125 The optimized structures were hydrated and subjected to molecular dynamics simulations using state-of-the-art simulation conditions (see Methods for detailed parametrization). 126 127 Colibactin appeared as a rather flexible molecule with an average end-to-end distance of around 13 Å (Extended Data Fig. 5B). The development of a putative model (see Methods) 128 shows a very stable binding of colibactin to the minor groove (Fig. 2C). It displays excellent 129 130 van der Waals contacts with all the groove's walls, with the cyclopropane rings pointing towards the adenines on opposite strands (Fig. 2D). Based on the equilibrium trajectory, we 131

132 determined that the number of base-pairs involved in the binding could fluctuate between 4 and 5, depending on the orientation of the cyclopropane ring and the carbon that alkylates 133 the N3 of the adenines (Fig. 2D, enlargement). In all cases, colibactin fitted perfectly into the 134 narrow minor groove of the targeted sequences and adopted a spatial arrangement with a 135 distance of around 4-5.5 Å between N3 and the center of mass of the cyclopropane, which 136 137 would facilitate N3 alkylation. These results and further analysis of sequence enrichment (Extended Data Fig. 5C) indicate that apart from the presence of adenines in a distance of 3-4 138 139 base-pairs on opposite strands and global DNA shape parameters (see above), the presence 140 of central A/T dinucleotides also influenced colibactin binding.

Using our approach depicted in Extended Data Fig. 1C, we investigated whether a specific 141 mutational signature associated with this sequence exists in cancers that have been 142 experimentally connected to pks+ E. coli infection <sup>6,32</sup>. Thus, we used whole-exome 143 sequencing data from CRC samples <sup>33</sup> (WXS project, n=619) and various databases (TCGA 144 145 project, https://www.cancer.gov/tcga, 10,224 tumor cases in 24 cancer types, n=553 for CRCs) to test whether somatic mutations are specifically enriched in the sequences 146 147 constituting the CDM (Methods). We found that mutation rates in AAWWTT motifs were 148 enriched as compared with all other WWWWW motifs in CRCs in both analyzed data sets 149 (Methods, Fig. 3A, Extended Data Fig. 6A). We also observed enrichment at AAWWTT motifs 150 in some other cancers (e.g., stomach cancer and uterine corpus endometroid cancer) and in 151 cancers containing polymerase epsilon (POLE) mutations (Extended Data Fig. 6A). A thorough inspection of mutated motifs in these cancers, however, indicated that at least some of this 152 153 enrichment could be attributed to other mutational processes also affecting the same motif (Extended Data Fig. 7A, B). 154

We validated the findings from the WXS data in another cohort of CRC samples, assessed by 155 whole genome sequencing (WGS)<sup>34</sup>, and analyzed the enrichment of mutations in CDM-156 related sequences for 200 tumors, including 184 microsatellite-stable (MSS), 3 POLE-157 mutated and 13 microsatellite-instable (MSI) cases. This analysis was run in a similar manner 158 159 to the WXS data analysis; however, each sample was analyzed separately instead of pooling 160 the samples into sub-cohorts, so as to identify motif enrichment and mutational loads in individual samples. We found significant (Mann-Whitney-U test, p<0.05, FDR <20%) 161 162 enrichment of mutations at colibactin-specific pentanucleotide sequences as compared with other sequences with the same length and A/T content in 3/3 POLE-mutated samples and 163

164 48/184 (26%) MSS cases, but not in MSI cases (Fig. 3B). We further found similar 165 enrichments in MSS samples by probing for colibactin-specific pentanucleotide 166 (AAATT/AAAAT) and hexanucleotide sequences (Fig. 3B and Extended Data Fig. 8A), whereby 167 pentanucleotide sequences showed less enrichment in MSI cases. The median number of 168 mutations in MSS samples at the pentanucleotide motifs was 963 (range: 66-11876), 169 corresponding to a median proportion of 6.7% (range: 3.9-44.7%).

We next compared mutations at the AAWWTT motifs to established mutational signatures 170 (Extended Data Fig. 1C) that have been previously deduced in pan-cancer analyses <sup>35</sup> and 171 normal colorectal epithelial cells  $^{3}$ . To this end, we derived a theoretical single base 172 substitution (SBS) signature based on the frequency of trinucleotide sequences, assuming a 173 174 uniform distribution of T>A, T>C and T>G mutations across the motifs (Methods, Fig. 3C). 175 Our comparison showed the highest similarity (cosine correlation coefficient 0.8) with SBSA, an SBS previously identified in healthy colonic mucosa  $^{3}$ . This signature is mainly 176 characterized by T>C mutations at ATA, ATT and TTT trinucleotides and T>G mutation at TTT 177 <sup>3</sup>. A total of 18 out of 184 MSS CRC cases (9.8%) showed SBSA contributions to the total 178 179 mutational load that was greater than 5%. SBS exhibiting preferential nucleotide changes at the same trinucleotides (SBS41, SBS57, SBS28) showed similarities as well (Fig. 3C). The 180 181 estimated contributions of SBSA to the total mutation load in the WGS of the CRC cohort correlated significantly with the respective proportion of mutated CDMs (Fig. 3D). This 182 183 finding suggests the SBSA trinucleotide signature could be a surrogate of the colibactinspecific mutational signature. Interestingly, we observed that the proportion of mutations at 184 185 CDMs was particularly enriched in sigmorectal tumor sites (Figure 3E, Kruskal-Wallis test, p <0.001). 186

To assess whether mutations in CDM could contribute to tumorigenesis by affecting pancancer driver genes, we searched for such mutations in the TCGA dataset. Intriguingly, CDMassociated driver gene mutations were identified in 18/286 (6.9%) of colon adenocarcinoma cases and in 8/89 (9.0%) of rectal adenocarcinoma cases, amongst which 10 out of 28 mutations were identified in the *APC* gene (Fig. 3F). Interestingly, we also observed frequent mutations in driver genes of other cancer types, of which the respective tissue might have been exposed to *E. coli* infection (Extended Data Fig. 8B and Supplementary Table 2).

We hypothesized that determining the exact ends of colibactin-induced DSBs (Extended DataFig. 1D) could provide additional information on the resulting mutational signature. Thus, we

devised a method for detecting the authentic 5' ends of a DSB in both forward and reverse 196 directions (Fig. 4A, Extended Data Fig. 9A), which could discriminate between blunt-ended 197 DSBs and DSBs with 5' or 3' overhangs (Extended Data Fig. 9A). The approach was validated 198 using published BLISS data from a cell line that contained an inducible endogenous AsiSI 199 restriction enzyme <sup>36</sup> (Extended Data Fig. 9B, Methods). To identify colibactin-specific break 200 201 patterns, we assessed the positional distribution of identified DSB ends at three CDM-related 202 sequences and three control sequences of the same length and similar AT content (Fig. 4B). All sequences revealed similar profiles of DSB ends in control pks- E. coli-infected cells and 203 204 untreated cells (Fig. 4B, C). Conversely, in *pks+ E. coli*-infected cells, we noticed a distinct 205 positional increase of DSB ends in colibactin-specific AAWWTT motifs as compared to control sequences (Fig. 4B, C). Interestingly, the observed colibactin-specific positional end point 206 207 increases suggest an action of colibactin resulting in 5' overhangs of 2 nucleotides, with a 208 preference of end points downstream of positions 2 and 5 in the hexamer motif of the 209 forward and reverse strands, respectively.

210 Next, we explored (Extended Data Fig. 1D) a possible relationship between the positions of 211 colibactin-induced damage and mutational changes in MSS CRC patients of the WGS cohort <sup>34</sup>. By comparing the 20 tumors with the highest vs. lowest mutational burden at CDMs, we 212 213 identified the highest number of mutations at positions 2 and 5 in the CDM-related sequences AAAATT and AATATT (Fig. 4D). For the AAATTT sequence motif, we found a 214 215 broader distribution of mutated nucleotide positions, compatible with the relatively broader distribution of DSB sites (Fig. 4D). The mutational changes at positions 2 and 5 nicely 216 217 correlated with the experimentally-deduced downstream (5'>3') break ends, which could be the result of single-strand breaks downstream (5'>3') to the damaged bases (Fig. 4E). 218 219 Mutational changes at those positions included mainly T:A>C:G. In agreement with previous information on SBSA<sup>3</sup>, we found that T:A>C:G mutations at CDM had a higher average allelic 220 221 frequency than similar changes at other motifs (two-sided Mann-Whitney-U test, p < 0.001) or than mutations attributable to the age-dependent signature SBS1<sup>35</sup> (two-sided Mann-222 223 Whitney-U test, p <0.001), indicating a mutational process happening early in carcinogenesis 224 rather than an ongoing process.

Using an unbiased global approach to assess the specificity of colibactin-induced DSB formation, we identified the motif AAWWTT as the preferred target. Most but not all sequences constituting the CDM exhibited extreme physical values of the DNA duplex,

characterized by a very narrow minor groove width, a highly negative electrostatic potential 228 and a pronounced stiffness. Even though a definitive understanding of the binding 229 230 conformation requires further support from 3D structural analysis, the 3D model provides valuable insight. In concert with our experimental data, this model suggests that colibactin, 231 232 using its two cyclopropane groups, selectively targets two adenines on opposite strands at a distance of ~4 nt in the motif AAWWTT. This model is in perfect agreement with the 233 234 observed DSB ends generated by colibactin's action and the position of mutations in the CDMs of the relevant cancer genomes. Accordingly, our data illuminate the mechanism of 235 236 DSB formation by colibactin and provide a firm link to the presence of a distinct mutational 237 signature in human cancer genomes, most specifically in colorectal carcinoma.

Recent mechanistic data suggest that, following the formation of N3 adenine adducts <sup>19</sup>, 238 colibactin's DNA cross-links undergo depurination via beta-elimination<sup>21</sup>. This may involve 239 the action of glycosylases <sup>37</sup>. Our exact DSB end point determination by sBLISS suggests that 240 241 cleavage occurs immediately downstream (5'>3') of positions 2 and 5 of the putativelyalkylated adenines on the opposing strands. Thus, cleavage results in staggered ends with 5' 242 243 overhangs of two nucleotides each (Extended Data Fig. 1E). Most intriguingly, the identified cleavage positions next to the alkylated adenines perfectly match the location of CDM-244 245 associated mutations in CRCs. The ultimate manifestation of mutations in the damaged lesions has been suggested to involve the action of error-prone translesion polymerases in 246 opposition to the damaged bases or apurinic sites <sup>38</sup>. Alternatively, other host DNA repair 247 mechanisms might act, such as nucleotide excision of alkylated adenines <sup>39</sup>. This could lead 248 249 to DSBs, resection of break ends, or complete repair, amongst others; however, details of the mutational process requires further investigation. 250

251 Most strikingly, we were able to demonstrate the enrichment of mutations in the CDMs of cancer genomes. Owing to the unique features of these mutations, we propose that this 252 253 provides definitive evidence of a bacterial infection signature in human cancer. Further, our observation of colibactin-related driver mutations in CRCs is suggestive of a potential causal 254 relationship. Typically for CRC, the initial driver mutations cause an intensification of Wnt 255 signaling, consistent with mutations in  $APC^{40}$ . In fact, our search for driver mutations in the 256 257 CDM of cancers revealed enrichment in such driver genes, particularly APC - thus 258 corroborating the causative role of pks+ *E. coli* in colon cancer.

- 259 An intriguing speculation is that, besides pks+ E.coli, other bacteria, such as Klebsiella
- 260 *pneumoniae* potentially harboring the same genotoxin island<sup>2</sup>, might be involved in human
- 261 carcinogenesis as well. While current data only provide weak evidence for a broader impact
- 262 of the colibactin in other cancers, this possibility remains as an important perspective for
- 263 future investigation.
- 264

### 265 References

- Nougayrede, J.P., et al. Escherichia coli induces DNA double-strand breaks in eukaryotic cells.
   *Science* 313, 848-851 (2006).
   Putze, J., et al. Genetic structure and distribution of the colibactin genomic island among members of the family Enterobacteriaceae. Infect Immun 77, 4696-4703 (2009).
- Lee-Six, H., *et al.* The landscape of somatic mutation in normal colorectal epithelial cells.
   *Nature* 574, 532-537 (2019).
- Watanabe, T., Tada, M., Nagai, H., Sasaki, S. & Nakao, M. Helicobacter pylori infection
   induces gastric cancer in mongolian gerbils. *Gastroenterology* **115**, 642-648 (1998).
- Castellarin, M., *et al.* Fusobacterium nucleatum infection is prevalent in human colorectal
   carcinoma. *Genome Res* 22, 299-306 (2012).
- Arthur, J.C., *et al.* Intestinal inflammation targets cancer-inducing activity of the microbiota.
   *Science* 338, 120-123 (2012).
- 278 7. Cougnoux, A., *et al.* Bacterial genotoxin colibactin promotes colon tumour growth by
  279 inducing a senescence-associated secretory phenotype. *Gut* 63, 1932-1942 (2014).
- 280 8. Bleich, R.M. & Arthur, J.C. Revealing a microbial carcinogen. *Science* **363**, 689-690 (2019).

Hibbing, M.E., Fuqua, C., Parsek, M.R. & Peterson, S.B. Bacterial competition: surviving and
thriving in the microbial jungle. *Nat Rev Microbiol* 8, 15-25 (2010).

- Takahashi, I., *et al.* Duocarmycin A, a new antitumor antibiotic from Streptomyces. *J Antibiot*(*Tokyo*) 41, 1915-1917 (1988).
- 11. Igarashi, Y., et al. Yatakemycin, a novel antifungal antibiotic produced by Streptomyces sp.
  TP-A0356. J Antibiot (Tokyo) 56, 107-113 (2003).
- Arcamone, F., Penco, S., Orezzi, P., Nicolella, V. & Pirelli, A. STRUCTURE AND SYNTHESIS OF
   DISTAMYCIN A. *Nature* 203, 1064-1065 (1964).
- Finlay, A., Hochstein, F., Sobin, B. & Murphy, F. Netropsin, a new antibiotic produced by a
  Streptomyces. *Journal of the American Chemical Society* **73**, 341-343 (1951).
- 29114.Boger, D.L. & Johnson, D.S. CC-1065 and the duocarmycins: unraveling the keys to a new292class of naturally derived DNA alkylating agents. *Proc Natl Acad Sci U S A* **92**, 3642-3649293(1995).
- 29415.Zha, L., et al. Colibactin assembly line enzymes use S-adenosylmethionine to build a295cyclopropane ring. Nat Chem Biol 13, 1063-1065 (2017).
- Healy, A.R., Vizcaino, M.I., Crawford, J.M. & Herzon, S.B. Convergent and Modular Synthesis
  of Candidate Precolibactins. Structural Revision of Precolibactin A. *J Am Chem Soc* 138, 54265432 (2016).
- 299 17. Zha, L., Wilson, M.R., Brotherton, C.A. & Balskus, E.P. Characterization of Polyketide Synthase
  300 Machinery from the pks Island Facilitates Isolation of a Candidate Precolibactin. ACS Chem
  301 Biol 11, 1287-1295 (2016).
- 302 18. Vizcaino, M.I. & Crawford, J.M. The colibactin warhead crosslinks DNA. *Nat Chem* 7, 411-417
  303 (2015).
- Wilson, M.R., *et al.* The human gut bacterial genotoxin colibactin alkylates DNA. *Science* **363**(2019).

| 306 | 20. | Xue, M., et al. Structure elucidation of colibactin and its DNA cross-links. Science <b>365</b> (2019).     |
|-----|-----|---|
| 307 | 21. | Xue, M., Wernke, K. & Herzon, S.B. Depurination of colibactin-derived interstrand cross-links.              |
| 308 |     | Biochemistry (2020).  |
| 309 | 22. | Yan, W.X., et al. BLISS is a versatile and quantitative method for genome-wide profiling of                 |
| 310 |     | DNA double-strand breaks. Nat Commun 8, 15058 (2017).   |
| 311 | 23. | Canela, A., et al. Genome Organization Drives Chromosome Fragility. Cell 170, 507-521.e518                  |
| 312 |     | (2017).   |
| 313 | 24. | Yang, F., Kemp, C.J. & Henikoff, S. Anthracyclines induce double-strand DNA breaks at active                |
| 314 |     | gene promoters. <i>Mutat Res</i> <b>773</b> , 9-15 (2015).  |
| 315 | 25. | Sheffield, N.C. & Bock, C. LOLA: enrichment analysis for genomic region sets and regulatory                 |
| 316 |     | elements in R and Bioconductor. Bioinformatics 32, 587-589 (2016).  |
| 317 | 26. | Bailey, T.L. DREME: motif discovery in transcription factor ChIP-seq data. <i>Bioinformatics</i> 27,        |
| 318 |     | 1653-1659 (2011).   |
| 319 | 27. | Tubbs, A., et al. Dual Roles of Poly(dA:dT) Tracts in Replication Initiation and Fork Collapse.             |
| 320 |     | <i>Cell</i> <b>174</b> , 1127-1142 e1119 (2018).  |
| 321 | 28. | Nelson, H.C., Finch, J.T., Luisi, B.F. & Klug, A. The structure of an oligo(dA).oligo(dT) tract and         |
| 322 |     | its biological implications. Nature <b>330</b> , 221-226 (1987).  |
| 323 | 29. | Yuan, G.C., et al. Genome-scale identification of nucleosome positions in S. cerevisiae.                    |
| 324 |     | Science <b>309</b> , 626-630 (2005).  |
| 325 | 30. | Tse, W.C. & Boger, D.L. Sequence-selective DNA recognition: natural products and nature's                   |
| 326 |     | lessons. Chem Biol 11, 1607-1617 (2004).  |
| 327 | 31. | Drsata, T., et al. Mechanical properties of symmetric and asymmetric DNA A-tracts:                          |
| 328 |     | implications for looping and nucleosome positioning. <i>Nucleic Acids Res</i> <b>42</b> , 7383-7394 (2014). |
| 329 | 32. | Buc, E., et al. High prevalence of mucosa-associated E. coli producing cyclomodulin and                     |
| 330 |     | genotoxin in colon cancer. <i>PLoS One</i> <b>8</b> , e56964 (2013).  |
| 331 | 33. | Giannakis, M., et al. Genomic Correlates of Immune-Cell Infiltrates in Colorectal Carcinoma.                |
| 332 |     | <i>Cell Rep</i> <b>15</b> , 857-865 (2016).   |
| 333 | 34. | Katainen, R., et al. CTCF/cohesin-binding sites are frequently mutated in cancer. Nat Genet                 |
| 334 |     | <b>47</b> , 818-821 (2015).   |
| 335 | 35. | Alexandrov, L.B., et al. The repertoire of mutational signatures in human cancer. Nature 578,               |
| 336 |     | 94-101 (2020).  |
| 337 | 36. | Iannelli, F., et al. A damaged genome's transcriptional landscape through multilayered                      |
| 338 |     | expression profiling around in situ-mapped DNA double-strand breaks. Nat Commun 8,                          |
| 339 |     | 15656 (2017).   |
| 340 | 37. | Iyama, T. & Wilson, D.M., 3rd. DNA repair mechanisms in dividing and non-dividing cells. DNA                |
| 341 |     | Repair (Amst) <b>12</b> , 620-636 (2013).   |
| 342 | 38. | Choi, J.Y., Lim, S., Kim, E.J., Jo, A. & Guengerich, F.P. Translesion synthesis across abasic               |
| 343 |     | lesions by human B-family and Y-family DNA polymerases alpha, delta, eta, iota, kappa, and                  |
| 344 |     | REV1. J Mol Biol <b>404</b> , 34-44 (2010).   |
| 345 | 39. | Martin, L.P., Hamilton, T.C. & Schilder, R.J. Platinum resistance: the role of DNA repair                   |
| 346 |     | pathways. <i>Clin Cancer Res</i> <b>14</b> , 1291-1295 (2008).  |
| 347 | 40. | Morin, P.J., et al. Activation of beta-catenin-Tcf signaling in colon cancer by mutations in                |
| 348 |     | beta-catenin or APC. Science <b>275</b> , 1/8/-1/90 (1997).   |
|     |     |   |

### 350 Acknowledgements

351 The results shown here are in whole or part based upon data generated by the TCGA 352 Research Network: https://www.cancer.gov/tcga. The authors would like to thank Silvano Garnerone from N.C. laboratory for processing raw sBLISS data, Pablo D. Dans from 353 354 Biophysical Chemistry Lab, Department of Biological Science (CENUR Litoral Norte), UdelaR, 355 UY, for constructive discussions about the theoretical model of colibactin, Ulrich Dobrindt from University of Münster for providing E. coli strains and Kfir Lapid and Rike Zietlow for 356 357 editing the manuscript. We also acknowledge the computational resources provided by the CSC-IT Center for Science, Finland. 358

359

### 360 Funding

361 P.J.D.K. was supported by IMPRS-IDI graduate school. B.A.M.B was supported by a Rubicon 362 fellowship from the Netherlands Organisation for Scientific Research (NWO). M.O. is an ICREA (Institució Catalana de Recerca i Estudis Avancats) academia researcher. This work was 363 364 supported by the German Research Foundation (DFG, grant number ME705/18-1) to T.F.M., by the Spanish Ministry of Science (grants BIO2015-64802-R, BFU2014-61670-EXP and 365 366 BFU2014-52864-R), the Catalan Government (grant 2014-SGR), the Instituto de Salud Carlos III-Instituto Nacional de Bioinformática (ISCIII PT 13/000/0030), the Biomolecular and 367 368 Bioinformatics Resources Platform and the EU Horizon 2020 research and innovation program (grants Elixir-Excelerate 676559 and BioExcel2:823830), and the MINECO Severo 369 370 Ochoa Award of Excellence to the IRB Barcelona, and by the Ragnar Söderberg Foundation, the Swedish Foundation for Strategic Research (BD15-0095) and the Strategic Research 371 372 Programme in Cancer (StratCan) at Karolinska Institutet to N.C. L.A.A. was supported by 373 grants from the Academy of Finland (Finnish Center of Excellence No. 312041, Academy 374 Professor grants 319083 and 320149, iCAN Flagship 320185), Jane and Aatos Erkko Foundation, Sigrid Juselius Foundation, Helsinki Institute of Life Science, and the Cancer 375 376 Society of Finland.

377

### 378 Author contributions

P.J.D.-K., H.B. and T.F.M. conceived the project and designed experiments, N.C. and B.A.M.B.
designed sBLISS experiments, P.J.D.-K. and B.A.M.B. performed sBLISS experiments, A.I.
performed E.coli infection and immunostaining experiments. P.J.D.-K. and A.I. prepared
samples for sBLISS. Bioinformatics analysis were performed by P.J.D.-K. and H.B. Theoretical
model of colibactin was built by F.B. and M.O. R.K., T.C. and L.A.A. provided and analyzed
WGS CRC data. The manuscript was written by P.J.D.-K., H.B. and T.F.M.

385

### 386 Competing interests

387 The authors declare no competing interests

388

### 389 Data and materials availability

Input FASTA files have been submitted to GEO under accession GSE145594 and analysis
 scripts are available from https://github.com/MPIIB-Department-TFMeyer/Dziubanska Kusibab\_et\_al.\_Colibactin. All other data is available in the main text or supplementary
 materials.

### 394 Figure Legends

395

### 396 Fig.1. Colibactin preferentially damages DNA in specific AT-rich motifs

397 (A) Illustration of the sBLISS protocol for the identification of colibactin-induced DSBs. (B) We 398 employed DREME to analyze motif enrichment around DSBs in pks+ vs. pks- E. coli-infected 399 cells. The sequence AAWWTT was identified as the preferred CDM among the top three enriched motifs in four independent replicates. (C) Relative enrichment of hexanucleotide 400 sequences in close proximity to the DSB positions (±7 nt) upon different treatments was 401 determined as the ratio of counts of a given sequence between two conditions. (D) Bars 402 403 represent the mean enrichment of the top 50 DSB-associated hexanucleotide motifs (log2 404 ratios of DSBs at each motif comparing two conditions) in pks+ vs. pks- E. coli-infected cells 405 (top) or pks- E. coli-infected cells vs. nontreated (NT) cells (bottom), n=4 independent experiments. Insets - results for all hexanucleotides. Error bars denote the 95% confidence 406 407 interval (CI) around the mean log2 ratios.

408

### 409 Fig. 2. Colibactin-specific binding motif correspond to extreme DNA shape parameter

### 410 values and electrostatic potential

(A) Relations between hexanucleotide sequence enrichments (log2 ratio standardized to 1;
mean of 4 independent experiments) in *pks+* vs. *pks- E.coli*-infected cells and values of
predicted DNA shape parameters for the central nucleotides of those motifs. Error bars
denote the 95% CI around mean scaled log2 ratios for each hexanucleotide sequence, n=4.

(B) 3D principal component analysis (PCA) visualization of the first three principal
components from predicted DNA shape values for all positions in all possible 6 nt motifs. The
motifs that match AAWWTT and/or show strong enrichment are highlighted: red –
significant AAWWTT sequences; blue – significant non-AAWWTT sequences; grey – other
sequences. Statistical significance is determined by an upper 95% CI limit of mean log2 ratio
> 0.2, n=4 independent experiments.

421 (C-D) Graphical representations of the theoretical docking of the predicted colibactin 422 structure into its preferred sequence motif (central sequence AAATTT), showing the 423 insertion of colibactin into the minor groove, with the double-stranded DNA as a surface (C)

or in atomic details (D). The enlargement shows a zoomed-in image of the proximity of the
cyclopropane to one of the N3 atoms of the adenine, highlighting the possibility to alkylate
the subsequent base-pair, depending on the carbon involved in the alkylation.

427

### 428 Fig. 3. Several cancers show enrichment of mutations in CDMs

429 Enrichment of SBS mutations at colibactin-specific AAATTT/AAAATT sequences in exomes of 619 CRC cases <sup>33</sup> (A) or at AAATT/AAAAT pentanucleotides in whole-genome data from 430 CRC<sup>34</sup> (B, n=200). (B) Top - differences in log2(mutations/base-pair) between colibactin-431 specific and all other WWWWW motifs. Middle - total mutation count at CDMs. Bottom -432 433 proportions of total mutations that overlap with CDMs in MSS, MSI, and POLE-mutated cases. (A,B) Asterisks - significant difference (two-sided Mann-Whitney-U test, p < 0.05 and 434 435 FDR < 20%) between CDMs and all WWWW(W) motifs. Numbers under asterisks – number of mutations overlapping AAWWTT / all mutations, [number of samples per cohort]. Error 436 437 bars =  $0 \pm 2MAD$  intervals of mutation rates (mutations/base-pair) of WWWW(W) sequences, excluding CDMs, after subtracting their mean. Dots - mutation rates for 438 439 AATTT/AAAAT or AAATTT/AAAATT after mean subtraction.. Crosses - means of the CDMs. Q1-Q4 cohorts defined by the quantiles of total SNV numbers. Outliers - samples with SNV 440 441 numbers significantly larger than the 95% CI in each cancer entity. POLE - polymerase epsilon mutated cases. (C) Theoretical trinucleotide change signature for uniform mutations at 442 443 AAWWTT motifs (top) and the matching signature SBSA (bottom).(D) Correlation between the estimated contribution of SBSA signature per sample in 184 MSS CRCs and the 444 445 proportion of total SNVs overlapping colibactin-specific pentanucleotides. (E) The proportion 446 of mutations in CDMs depending on the tumor location (n=184 MSS CRC, Kruskal-Wallis test, 447 p < 0.001). (F) APC exonic somatic mutations in CDM in CRC cases from TCGA and the 448 COSMIC database.

449

### 450 Fig. 4 Analysis of colibactin-induced DSB break ends and mutational patterns in CDMs

### 451 indicate damaged positions

(A) Schematic representation of the break end determination of an arbitrary DSB
configuration based on sBLISS data (see Fig. S5 for additional examples). (B) Mean
proportions of break ends by the position and direction of a motif in reads (red - forward;

green - reverse) for colibactin-specific (leftmost three columns) and control motifs 455 456 (rightmost three columns) in *pks+ E. coli*-infected cells, *pks- E. coli*-infected cells and control 457 conditions, n=4 independent experiments (dots). Negative values denote mean proportions on the reverse strand of the motif. Error bars: 95% CI around mean proportions. (C) The 458 mean differences in break end proportions (taken from Fig. 4B) between pks+ E. coli-infected 459 460 cells and pks- E. coli-infected cells, n=4. Negative values denote mean differences on the 461 reverse strand. Error bars: 95% CI around mean differences. Only increased break end proportions in *pks+ E. coli*-infected cells are shown for each position. (D) The positional 462 463 frequencies of nucleotide changes in CDMs found in the CRC WGS cohort. The 20 cases with 464 the highest and lowest proportion of mutations at AAWWTT are shown. (E) The model of colibactin-induced DSB and damaged bases as derived from the data shown in this figure. 465 466 Red-colored nucleotides denote damaged bases, vertical bars are single-strand breaks, which 467 were detected by sBLISS, and the broken line denotes the inter-strand crosslink caused by colibactin action. 468

469

- 470 Materials and Methods
- 471

### 472 Cell line, bacterial strains, E. coli infection and etoposide-treatment

Caco-2 cells (from ATCC<sup>®</sup> HTB-37<sup>™</sup>) were cultured at 37 °C in a humidified 5% CO<sub>2</sub> 473 474 atmosphere, in DMEM medium (Life Technologies, cat. number: 10938-025), supplemented with 20% FCS (Biochrom, cat. number: S0115). Contamination of Mycoplasma spp. in this 475 476 immortalized cell line was excluded using the Venor®GeM OneStep PCR kit (Minerva Biolabs<sup>®</sup>, cat. number: 11-8250). To infect Caco-2 cells, an overnight liquid culture of *E. coli* 477 478 strain M1/5 (streptomycin-resistant and colibactin-positive) and E. coli strain M1/5:: \Larbert clibactin-positive) and E 479 (streptomycin-resistant and colibactin-negative) was set up. Bacteria were inoculated in 5 ml 480 of Luria broth (LB) medium and incubated overnight at 37 °C in a shaking incubator. The 481 overnight inoculum was diluted 1:33 in infection medium (DMEM + 10% FCS + 10 mM HEPES 482 (Life Technologies, cat. number: 15630-056)) to obtain OD600=1 after 3 hr of incubation to give 1.5 x 10<sup>9</sup> bacteria/ml. The prepared bacteria inoculum was further diluted to reach MOI 483 20, added to Caco-2 cells seeded previously and incubated for 3 hr at 37 °C. Medium was 484 485 then aspirated and cells fixed according to the protocol for immunofluorescence or sBLISS.

For every biological replicate, positive (etoposide-treatment) and negative (no treatment) controls were included. Etoposide powder (Sigma Aldrich, cat. number: E1383) was diluted in DMSO to 50 mM and aliquots stored at -20 °C. Final drug dilutions to the concentration of 50  $\mu$ M were performed in pre-warmed infection medium prior to each treatment. Treatment was conducted for 3 hr at 37 °C, after which the medium was aspirated, and etoposidetreated cells were fixed in the same way as *E. coli*-infected cells.

492

### 493 Immunofluorescence staining

494 Caco-2 cells grown and infected on MatTek glass-bottom dishes were washed three times 495 with PBS (Life Technologies, cat. number: 14190-094) and fixed with 3.7% paraformaldehyde 496 (Sigma Aldrich, cat. number: P6148) for 1 h. The cells were kept overnight in blocking buffer 497 (3% BSA, Biomol, cat. number: 01400.100), 1% saponin (Sigma Aldrich, cat. number: 84510), 2% Triton X-100 (Carl Roth, cat. number: 3051.2) and 0.02% sodium azide (Sigma Aldrich, cat. 498 499 number: S2002). Blocking was followed by overnight incubation with yH2AX antibody (Phospho-Histone H2A.X (Ser139) Antibody, Cell Signaling, cat. number: 2577, 1:500 500 501 dilution) at 4 °C. The next day, the MatTek dishes were washed three times with blocking buffer followed by overnight incubation with secondary antibody (Dianova, cat. number: 502 503 711-035-152, 1:250 dilution) diluted in blocking buffer. Phalloidin 546 (Invitrogen, cat. number: A22283, 1:200 dilution) and Hoechst (Sigma, cat. number: H6024, 1:10000 dilution) 504 505 were added for staining actin filaments and DNA, respectively. The next day, cells were washed three times with blocking buffer and coverslipped using Vectashield<sup>®</sup> Antifade 506 507 Mounting Medium (Vector Laboratories, cat. number: H-1000). Images were acquired using 508 a Leica TCS SP-8 confocal microscope and processed using ImageJ.

509

### 510 sBLISS, an adaptation of the BLISS method

511 DSBs were identified using the suspension-cell BLISS (sBLISS) method <sup>41</sup>, which is an 512 adaptation of the previously published BLISS protocol <sup>42,43</sup>. In contrast to BLISS, where DSBs 513 are labeled in fixed cells immobilized on microscope slides, in sBLISS DSBs are labeled in fixed 514 cell suspensions. In brief, cells were treated/infected in culture dishes, after which they were 515 trypsinized, counted, centrifuged and resuspended in pre-warmed medium to obtain 10<sup>6</sup> 516 cells/ml. Then, cells were fixed by adding 16% PFA (Electron Microscopy Sciences, cat.

517 number: 15710) to reach a final concentration of 4%. After 10 minutes, 2 M glycine 518 (Molecular Dimensions, cat. #MD2-100-105) was added to a final concentration of 125 mM 519 in order to block unreacted aldehydes. This was followed by two 5 min incubations, first at 520 room temperature and then on ice, followed by two washes in ice-cold PBS. Cross-linked 521 cells were stored in PBS at 4 °C until further processing.

522 Next, the BLISS template was prepared. This includes (1) Cell lysis in 10 mM Tris-HCl, 10 mM 523 NaCl, 1 mM EDTA, and 0.2% Triton X-100 (pH 8) buffer, followed by lysis in buffer containing 524 10 mM Tris-HCl, 150 mM NaCl, 1 mM EDTA, and 0.3% SDS (pH 8); (2) DSB blunting with NEB's Quick Blunting Kit (NEB, cat. number: E1201); (3) In situ BLISS adapter ligation using T4 DNA 525 Ligase (ThermoFisher Scientific, cat. number: EL0011). Each BLISS adapter contained a T7 526 promoter sequence for IVT, the RA5 Illumina RNA adapter sequence, a random 8 nt long 527 528 sequence referred to as Unique Molecular Identifier (UMI) and an 8 nt long sample barcode; (4) Phenol:chloroform-based extraction of gDNA; (5) Fragmentation of isolated genomic DNA 529 (400-600bp) using BioRuptor Plus (Diagenode). Obtained BLISS templates were stored at -20 530 °C. 531

The final step of the BLISS protocol was in vitro transcription (IVT) followed by NGS library 532 533 preparation. First, 100 ng of purified, sonicated and differentially-barcoded BLISS template of 1) etoposide-treated and non-treated cells, or 2) cells infected with pks+ E. coli or with pks-534 535 E. coli were pooled into one reaction, respectively. IVT was performed using MEGAscript T7 Transcription Kit (ThermoFisher, cat. number: AMB13345) for 14 hr at 37 °C in the presence 536 537 of RiboSafe RNAse Inhibitor (Bioline, cat. number BIO-65028). Next, gDNA was removed using DNase I (ThermoFisher, cat. number: AM2222), and the remaining RNA was purified 538 539 with Agencourt RNAClean XP beads (Beckman Coulter). The Illumina RA3 adapter sequence was ligated to the purified RNA using T4 RNA Ligase 2 (NEB, cat. number: M0242) for 2 hr at 540 541 25 °C and reverse transcription was performed with Reverse Transcription Primer (Illumina sequence) using SuperScript IV Reverse Transcriptase (ThermoFisher, cat. number: 542 543 18090050) for 50 minutes at 50 °C. This was followed by heat inactivation for 10 minutes at 544 80 °C. Finally, libraries were amplified with NEBNext High-Fidelity 2x PCR Master Mix (NEB, 545 cat. number: M0541), the RP1 common primer and a uniquely selected index primer. 12 PCR 546 cycles were conducted, and after that, libraries were purified according to the two-sided 547 AMPure XP bead purification protocol (Beckman Coulter). Profiles of the libraries were

quantified on a BioAnalyzer High Sensitivity DNA chip. Libraries were sequenced as singleend (1x75) reads on the NextSeq platform.

550

### 551 Pre-processing of sequencing data

Raw sequencing data were pre-processed as previously described <sup>42</sup>. In brief, only reads 552 which contained the expected prefix of UMI and sample barcode were kept using SAMtools 553 554 <sup>44</sup>. One mismatch in the barcode sequence was allowed. Further, prefixes were trimmed, and the remaining sequences were aligned to the GRCh37/hg19 reference genome using BWA-555 MEM  $^{45}$ . Reads with mapping quality scores  $\leq$  30 and those that were determined as PCR 556 duplicates were removed. Finally, a BED file containing a list of unique DSB locations was 557 generated. DSBs which fell into ENCODE blacklist regions <sup>46</sup>, high coverage regions and low 558 mappability regions <sup>47</sup> were removed. Kept positions of DSBs were further used in 559 downstream analysis. For BLISS data of AsiSI-induced breaks <sup>48</sup> a similar procedure with 560 561 simplified filtering was applied.

562

### 563 Locus Overlap Analysis

To identify significant overlaps of DNA DSB with genomic region sets, we used LOLA<sup>11</sup>. We 564 first defined whole genome as a Universe Set, which was next divided into tiles of equal 565 566 lengths (1,000 nt). For each created tile, we next searched for overlaps with DSBs captured by sBLISS using the findOverlap() function. All tiles containing  $\geq$  10 breaks were used as a 567 568 Query Set. The runLOLA() function was executed with LOLA Core databases (reduced by 569 Tissue clustered DNase hypersensitive sites) as well as LOLA Extended databases and custom 570 database containing non-B-DNA regions (https://nonbabcc.ncifcrf.gov/apps/site/references). Fisher's exact test was used with a FDR  $\leq$  5%. 571

572

### 573 Pattern identification around breakpoints

We screened the +/-10 bp contexts of DSB positions produced by sBLISS using DREME<sup>49</sup>, comparing all breakpoints from pks+ E. *coli*-infected cells to pks- E. *coli*-infected cells from the same experiment, allowing for a maximum of 10 motifs per experiment with a minimum E value of less than 0.05. A similar analysis was performed for etoposide vs. untreated

578 control cells but allowing for up to 50 motifs per experiment. Most of the top motifs 579 identified for pks+ E.coli- vs. pks- E.coli-infected cells were of 7 nt length; therefore the 580 subsequent computation of enrichment log ratios for all possible oligonucleotides was done on windows of DSB breakpoints +/-7 bp and for oligonucleotide lengths of 5-7 nt. We 581 582 computed the proportion of all DSB breakpoints in each sample that contained a given 583 pattern and determined the log2 ratio of the proportion of each pattern in the test 584 compared to control conditions of the same experiment. We observed that several of the 585 patterns enriched only in specific experiments were very similar to the library barcode from 586 the same experiments, indicating a priming activity of unbound sBLISS linkers. We therefore 587 excluded all patterns matching library-specific barcodes with at most 1 mismatch from further analysis. Enriched heptamers were, in agreement with the DREME analysis, mostly 588 589 the top enriched hexamers with different flanking nucleotides and are therefore not 590 reported here. To visualize the proportion of each nucleotide per position the package seqLogo was used 50,51 591

592

### 593 **DNA Shape predictions**

594 DNA structures can be described in terms of base-pair and base-step parameters that consist 595 of three translational and rotational movements between the bases or the base pairs, respectively. At the base-pair step level, DNA deformability along these six directions has 596 been described by the associated stiffness matrix <sup>52</sup>. From the ensemble of molecular 597 dynamics (MD) simulations considering the tetramer environment using the newly refined 598 599 parmbsc1 force field, we retrieved the 6x6 matrix describing the deformability of the helical 600 parameters for each possible DNA tetramer. Pure stiffness constants corresponding to the 601 six base-pair step parameters (shift, slide, rise, tilt, roll and twist) were extracted from the 602 diagonal of the matrix, and the total stiffness ( $K_{\rm L}$  tot) was obtained as a product of these six 603 constants and used as an estimate of the flexibility of each base pair step in a tetramer. For 604 predictions of minor groove width (MGW), propeller twist (ProT), electrostatic potential (EP), 605 helical twist (HeIT) and roll (Roll) the getShape function from 'DNAshapeR' package was used <sup>53</sup>. Input FASTA files, containing sequences in close proximity to identified DSB (±5 nt), were 606 607 extracted with a custom Python script. Central values of DNA shape parameters and EP were 608 derived as the mean values of nucleotide at positions 3 and 4 of the motifs for MGW, ProT 609 and EP; for stiffness, HeIT and Roll the mean of the base-pair values is calculated for steps 2

to 4 in each hexamer. For the multivariate analysis of those parameters we used all values along the full 6nt of hexanucleotide sequences predicted in the same way. The interaction potential (electrostatic and van der Waals) of Na<sup>+</sup> probes with DNA duplexes was determined using a linear approximation to the Poisson-Boltzmann equation and dielectric constant for the DNA as implemented in the CMIP program <sup>54</sup>.

615

### 616 Model and Molecular dynamics set up

The 3D structure and protonation state of colibactin were built starting from the smile 617 (https://pubchem.ncbi.nlm.nih.gov/compound/138805674#section=InChI) 618 using 619 MarvinSketch (MarvinSketch, version 6.2.2, calculation module developed by ChemAxon, 620 http://www.chemaxon.com/products/marvin/marvinsketch/). The geometry of the model and the partial atomic charges were assigned to the structure with General Amber Force 621 Field (GAFF) <sup>55</sup>. Parameters and topology files were prepared with Acpype <sup>56</sup>. The colibactin 622 was then simulated in explicit solvent at 298 K (see below for details) for 250 ns, and along 623 624 the simulation the distance between the cyclopropanes was monitored (see Extended Data Fig. 3B), to study their orientation and the overall length of the free colibactin. Using 625 HADDOCK 2.4<sup>57</sup>, we then built the DNA-colibactin complex. For the docking, we selected a 626 representative structure of free colibactin along the MD simulation, with an average distance 627 among the cyclopropanes (red line, Extended Data Fig. 3B) and an equilibrated structure of 628 629 the DNA (sequence CGAAATTTCG). After the initial docking, which positioned the molecule correctly along the minor groove of the DNA, we then manually rotated the molecule slightly 630 631 to improve the orientation of the cyclopropanes towards the N3 of the closest adenine using 632 PYMOL (The PyMOL Molecular Graphics System, Schrödinger, LLC (2018)). To check the 633 stability of this complex and to equilibrate its structure, the model was simulated (see details MD simulation below) and minimized in solution with positional restraints on the solute 634 using our well-established multi-step protocol <sup>58,59</sup>. The minimized structure was thermalized 635 to 298 K at NVT, and then simulated first applying harmonic restraints of 5 kcal/mol·Å2 on 636 637 the DNA structure and distance constraints between the cyclopropane and the N3 of the adenine (respectively 4 and 5 bases apart), each represented by a harmonic restraint of 2.5 638 639  $kcal/mol Å^2$ . To further check the stability of the complex, we then slowly removed the 640 constraints and ran 10 replicas of MD simulation of the complex for a total of 500 ns by 641 means of Molecular Dynamics simulations at NPT (P = 1 atm; T= 298K) and one MD

simulation 350 ns long. The first 10 ns of the simulations were considered as an equilibrationstep and were discarded for further analysis.

644 In each MD simulation (DNA, free colibactin and their complex, respectively), we placed the solute in the center of a truncated octahedral box of TIP3P water molecules <sup>60</sup>, neutralized 645 by K+ ions. In each simulation, the Berendsen algorithm <sup>61</sup> was used to control the 646 temperature and the pressure, with a coupling constant of 5 ps; and the SHAKE algorithm 647 was utilized to equilibrium the length of hydrogen atoms involved in the covalent bonds <sup>62</sup>. 648 649 Long-range electrostatic interactions were accounted for by using the Particle Mesh Ewald method (14) with standard defaults, and a real-space cut-off of 10 Å. For the DNA we used 650 the newly revised force field parmBSC1<sup>63</sup>. All simulations were carried out using AMBER 18 651  $^{\rm 64}$  , analyzed with CPPTRAJ  $^{\rm 65}$  and visualized using VMD 1.9.4  $^{\rm 66}.$ 652

653

### 654 Cancer somatic mutation data and mutation enrichment at colibactin-specific patterns

We obtained somatic variant data from the TCGA Unified Ensemble "MC3" Call Set <sup>67</sup> ("TCGA 655 pan-cancer dataset") and from the supplementary data of Giannakis et al. <sup>68</sup>. To test for 656 enrichment of mutations at any motif, we first identified positions of all hexanucleotide 657 658 motifs in the exonic portion of the genome. Somatic variants occurring at A or T bases were grouped in one of 6 classes (quartile 1-4, outlier or POLE mutated sample) depending on the 659 660 total SNV number and POLE mutation status of the corresponding tumor sample. We included all cases with an exonic mutation with predicted impact on the coded protein in 661 662 POLE and also those with similar mutations in the related gene POLD1 in the class of POLEmutated cases. We then computed the mutation rate for each hexanucleotide motif with 663 664 respect to the number of genomic bases covered in exonic regions for the same motif. As a 665 baseline, we established the mutation rates of all WWWWWW motifs and subtracted their 666 mean from the mutation rate of all other hexanucleotide motifs. We then tested for significance of the mutation rate at colibactin associated AAWWTT motifs (i.e., AAATTT and 667 668 AAAATT/AATTTT) compared to the remaining WWWWWW motifs using Mann-Whitney-U tests and computed the false discovery rate (FDR) using the method of Benjamini-Hochberg 669 670 <sup>69</sup>. We applied a similar procedure for pentanucleotide motifs AAATT/AATTT and AAAAT/ATTTT. Reads from WGS of colorectal cancers <sup>70</sup> (EGA database accession code 671 EGAS00001003010) were aligned to GRCh38 with BWA-MEM <sup>45</sup> and called using Mutect2 <sup>71</sup>. 672

All single nucleotide variant calls (PASSed by Mutect2) were used to determine the number 673 of mutations overlapping WWWWW pentanucleotides and WWWWWW hexanucleotides 674 and further analyzed in a similar way as for exome sequencing data on an individual sample 675 676 basis. For correlation of colibactin-specific mutation load with SBSA and clinical parameters, 677 we computed a colibactin mutation proportion. We considered colibactin mutations as all A/T > C/G mutations at 2nd and 5th position in the AAATTT, AATATT, AAAATT and AATTTTT 678 motifs. That number was divided by all other A/T > C/G mutations outside those motifs. 679 Somatic mutation and motif analyses with WGS CRC samples were performed with 680 BasePlayer<sup>72</sup>. Allelic proportions of somatic variants from rectal MSS WGS CRC classified as 681 colibactin associated (A/T > C/G mutations at 2nd and 5th position in the AAATTT, AATATT, 682 AAAATT and AATTTTT) were compared to those of all other somatic variants or age-related 683 684 signature SBS1 associated variants (C/G>T/A at CpG) of the same tumors using

685

### 686 Correlation with known trinucleotide signatures

Known trinucleotide SBS signatures were obtained from COSMIC Mutational Signatures v3<sup>73</sup> 687 Signatures from Lee-Six et al. <sup>74</sup> were generated using R scripts and trinucleotide frequencies 688 provided with this publication. We used the trinucleotide context of the three possible 689 690 sequences of the AAWWTT motif (AAATTT, AATTTT and AATATT) to determine trinucleotide frequencies in those motifs. We used the reverse complement of the trinucleotide if the 691 692 middle base was an A and extended the sequence by N (i.e., all possible nucleotides) for 693 trinucleotides overlapping the end of the motif by one base. For each possible trinucleotide, 694 we established the total number of occurrences in AAWWTT. We assumed a uniform 695 distribution of all possible mutations T>A/C/G nucleotide changes across all positions of the 696 motif, and then generated a trinucleotide change signature by computing the proportion of each combination of trinucleotide and nucleotide change in AAWWTT motifs. The cosine 697 698 correlation was computed between this AAWWTT signature and all known signatures in R. 699 Contributions to known trinucleotide signatures in cancer samples were estimated using the R package deconstructSigs <sup>75</sup>. 700

701

### 702 Driver gene analysis

703 Entity-specific driver genes were obtained from the Catalog of Cancer Genes (Cancer Genome Interpreter, https://www.cancergenomeinterpreter.org/genes)<sup>76</sup>. Pan-cancer driver 704 genes were obtained from a published analysis based on the TCGA cohort <sup>77</sup>. Somatic 705 variants from this database were classified as protein changing if the impact class was HIGH 706 707 or it was MODERATE, and the protein change was classified by PolyPhen or SIFT as probably damaging or deleterious, respectively. Somatic variants from samples of genome-wide 708 screens other than TCGA in COSMIC<sup>78</sup> were classified as protein changing if it had been 709 classified as pathogenic by FATHMM or if its consequence was a frameshift deletion or a 710 711 nonsense substitution.

712

### 713 Analysis of biases in AAWWTT mutation enrichment due to other signatures

We obtained frequencies of pentanucleotide changes, trinucleotide changes and 714 trinucleotide SBS signature contributions for PCAWG samples from the PCAWG project<sup>79</sup>. We 715 derived a theoretical AAWWTT pentanucleotide signature similar to the AAWWTT 716 trinucleotide signature above. We projected trinucleotide changes of known cancer-717 718 associated SBS signatures to trinucleotides, computing the weight of each pentanucleotide 719 change as the weight of the corresponding central trinucleotide change in the SBS signature 720 followed by normalization to a sum of one across all pentanucleotide changes. We computed the contribution of SBSA in all PCAWG samples using trinucleotide counts provided by 721 PCAWG and used the contributions of all other SBS signatures as provided by PCAWG. 722

723

### 724 Clinical data

Limited clinical data on tumor patients was obtained from a study by Katainen et al <sup>70</sup> that was reviewed and approved by the Ethics Committee of the Hospital District of Helsinki and Uusimaa (HUS). Signed informed consent or authorization from the Finnish National Supervisory Authority for Welfare and Health has been obtained for all the sample materials used.

730

### 731 Break end analysis

732 DSB positions were derived from sBLISS estimated breakpoints in a strand-specific manner as the middle of the positions of the first base in the mapped genomic fragment and the last 733 base in 5' direction (i.e. the direction of the linker). Genomic positions of oligonucleotide 734 patterns of interest were identified in GRCh37. Breakpoints falling with the pattern +/- 2bp 735 736 flanking regions were kept and the relative position of the breakpoint in the 5'-3' motif were 737 computed. Each strand-specific breakpoint was classified as mapping to the strand 738 containing the motif in forward or backward direction. Counts for each relative pattern position and motif direction were computed for each pattern and normalized by the total 739 number of breaks mapping to this pattern for each sample. Differences between normalized 740 741 counts were computed for each experiment between *pks+ E. coli* and *pks- E. coli* conditions.

742

### 743 Pattern enrichment around nucleosomes

Genomic positions of all hexanucleotide patterns with 6 A or T were identified in GRCh37. We used nucleosome positions identified in Gaffney at al <sup>80</sup> and determined the distance of each pattern to the nearest nucleosome center, keeping only those pattern sites within +/-1500bp around the center. Patterns were classified into three groups: a) AAWWTT motifs b) AAAAAA, ATATAT and TATATA which showed a distribution distinct from all other patterns around nucleosome centers and c) all remaining patterns.

750

### 751 Data analysis and visualization

All visualizations and statistical analyses were produced using R v3.4 <sup>81</sup>. For all boxplots, the top, center and bottom of the box correspond to the 75<sup>th</sup>, 50<sup>th</sup> (median) and 25<sup>th</sup> percentile of all data; the whiskers extend to the most extreme data points, which are no more than 1.5 times the length of the box (i.e., the inter-quartile range) beyond the box. All data points further away are shown as outliers.

757

### 758 Reporting Summary

Further information is available from the Life Sciences Reporting Summary linked to thisdocument.

### **References for Materials and Methods**

| 764 | 41.       | Gothe, H.J., et al. Spatial Chromosome Folding and Active Transcription Drive DNA Fragility                                   |
|-----|-----------|---|
| 765 |           | and Formation of Oncogenic MLL Translocations. <i>Mol Cell</i> <b>75</b> , 267-283.e212 (2019).                               |
| 766 | 42.       | Yan, W.X., et al. BLISS is a versatile and quantitative method for genome-wide profiling of                                   |
| 767 |           | DNA double-strand breaks. <i>Nat Commun <b>8</b>,</i> 15058 (2017).   |
| 768 | 43.       | Zhang, F., et al. Breaks Labeling in situ and sequencing (BLISS). Protocol Exchange DOI:                                      |
| 769 |           | 10.1038/protex.2017.018(2017).  |
| 770 | 44.       | Li, H., et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25, 2078-                                       |
| 771 |           | 2079 (2009).  |
| 772 | 45.       | Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.                                    |
| 773 |           | Bioinformatics <b>25</b> , 1754-1760 (2009).  |
| 774 | 46.       | An integrated encyclopedia of DNA elements in the human genome. Nature 489, 57-74   |
| 775 |           | (2012).   |
| 776 | 47.       | Speir, M.L., et al. The UCSC Genome Browser database: 2016 update. Nucleic Acids Res 44,                                      |
| 777 |           | D717-725 (2016).  |
| 778 | 48.       | Iannelli, F., et al. A damaged genome's transcriptional landscape through multilayered  |
| 779 |           | expression profiling around in situ-mapped DNA double-strand breaks. Nat Commun 8,  |
| 780 |           | 15656 (2017).   |
| 781 | 49.       | Bailey, T.L. DREME: motif discovery in transcription factor ChIP-seq data. Bioinformatics 27,                                 |
| 782 |           | 1653-1659 (2011).   |
| 783 | 50.       | Schneider, T.D. & Stephens, R.M. Sequence logos: a new way to display consensus   |
| 784 |           | sequences. Nucleic Acids Res 18, 6097-6100 (1990).  |
| 785 | 51.       | Bembom, O. seqLogo: Sequence logos for DNA sequence alignments. R package version   |
| 786 |           | 1.44.0. (2017).   |
| 787 | 52.       | Drsata, T., et al. Mechanical properties of symmetric and asymmetric DNA A-tracts:  |
| 788 |           | implications for looping and nucleosome positioning. Nucleic Acids Res 42, 7383-7394 (2014).                                  |
| 789 | 53.       | Chiu, T.P., et al. DNAshapeR: an R/Bioconductor package for DNA shape prediction and  |
| 790 |           | feature encoding. Bioinformatics 32, 1211-1213 (2016).  |
| 791 | 54.       | Gelpi, J.L., et al. Classical molecular interaction potentials: improved setup procedure in                                   |
| 792 |           | molecular dynamics simulations of proteins. <i>Proteins</i> <b>45</b> , 428-437 (2001).                                       |
| 793 | 55.       | Wang, J., Wolf, R.M., Caldwell, J.W., Kollman, P.A. & Case, D.A. Development and testing of a                                 |
| 794 |           | general amber force field. <i>J Comput Chem</i> <b>25</b> , 1157-1174 (2004).   |
| 795 | 56.       | Sousa da Silva, A.W. & Vranken, W.F. ACPYPE - AnteChamber PYthon Parser interfacE. BMC  |
| 796 |           | Res Notes 5, 367 (2012).  |
| 797 | 57.       | van Zundert, G.C.P., et al. The HADDOCK2.2 Web Server: User-Friendly Integrative Modeling                                     |
| 798 |           | of Biomolecular Complexes. J Mol Biol <b>428</b> , 720-725 (2016).  |
| 799 | 58.       | Dans, P.D., et al. Long-timescale dynamics of the Drew-Dickerson dodecamer. Nucleic Acids                                     |
| 800 |           | <i>Res</i> <b>44</b> , 4052-4066 (2016).  |
| 801 | 59.       | Perez, A., Luque, F.J. & Orozco, M. Dynamics of B-DNA on the microsecond time scale. J Am                                     |
| 802 | ~~        | Chem Soc <b>129</b> , 14/39-14/45 (2007).   |
| 803 | 60.       | Jorgensen, W.L., Chandrasekhar, J., Madura, J.D., Impey, R.W. & Klein, M.L. Comparison of                                     |
| 804 |           | simple potential functions for simulating liquid water. The Journal of chemical physics <b>79</b> ,                           |
| 805 | 64        | 926-935 (1983).   |
| 806 | 61.       | Berendsen, H.J., Postma, J.v., van Gunsteren, W.F., DiNola, A. & Haak, J.R. Molecular   |
| 807 |           | aynamics with coupling to an external bath. The Journal of chemical physics <b>81</b> , 3684-3690                             |
| 808 | <b>62</b> | (1984).<br>Diskovit k D. Chastilli C. B. Davas dava (k k Nassi i kitista sitisti i kitista sitisti i kitista sitisti i kitist |
| 809 | 62.       | Ryckaert, JP., Ciccotti, G. & Berendsen, H.J. Numerical integration of the cartesian equations                                |
| 810 |           | of motion of a system with constraints: molecular dynamics of n-alkanes. <i>Journal of</i>                                    |
| 811 |           | computational physics <b>23</b> , 327-341 (1977).   |

| 812<br>813 | 63. | Ivani, I., et al. Parmbsc1: a refined force field for DNA simulations. <i>Nat Methods</i> <b>13</b> , 55-58 (2016).  |
|------------|-----|--|
| 814        | 64. | Case, D., et al. AMBER 18. University of California, San Francisco (2018).   |
| 815<br>816 | 65. | Roe, D.R. & Cheatham, T.E., 3rd. PTRAJ and CPPTRAJ: Software for Processing and Analysis of Molecular Dynamics Trajectory Data. <i>J Chem Theory Comput</i> <b>9</b> , 3084-3095 (2013). |
| 817        | 66. | Humphrey, W., Dalke, A. & Schulten, K. VMD: visual molecular dynamics. J Mol Graph 14, 33-   |
| 818        |     | 38, 27-38 (1996).  |
| 819        | 67. | Ellrott, K., et al. Scalable Open Science Approach for Mutation Calling of Tumor Exomes Using  |
| 820        |     | Multiple Genomic Pipelines. Cell Syst 6, 271-281.e277 (2018).  |
| 821        | 68. | Giannakis, M., et al. Genomic Correlates of Immune-Cell Infiltrates in Colorectal Carcinoma.   |
| 822        |     | <i>Cell Rep</i> <b>15</b> , 857-865 (2016).  |
| 823        | 69. | Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful  |
| 824        |     | approach to multiple testing. Journal of the Royal statistical society: series B (Methodological)  |
| 825        |     | <b>57</b> , 289-300 (1995).  |
| 826        | 70. | Katainen, R., et al. CTCF/cohesin-binding sites are frequently mutated in cancer. Nat Genet  |
| 827        |     | <b>47</b> , 818-821 (2015).  |
| 828        | 71. | Cibulskis, K., et al. Sensitive detection of somatic point mutations in impure and   |
| 829        |     | heterogeneous cancer samples. Nat Biotechnol <b>31</b> , 213-219 (2013).   |
| 830        | 72. | Katainen, R., et al. Discovery of potential causative mutations in human coding and  |
| 831        |     | noncoding genome with the interactive software BasePlayer. Nat Protoc 13, 2580-2600  |
| 832        |     | (2018).  |
| 833        | 73. | Alexandrov, L.B., et al. The Repertoire of Mutational Signatures in Human Cancer. bioRxiv,   |
| 834        |     | 322859 (2018).   |
| 835        | 74. | Lee-Six, H., et al. The landscape of somatic mutation in normal colorectal epithelial cells.   |
| 836        |     | Nature <b>574</b> , 532-537 (2019).  |
| 837        | 75. | Rosenthal, R., McGranahan, N., Herrero, J., Taylor, B.S. & Swanton, C. DeconstructSigs:  |
| 838        |     | delineating mutational processes in single tumors distinguishes DNA repair deficiencies and  |
| 839        |     | patterns of carcinoma evolution. <i>Genome Biol</i> <b>17</b> , 31 (2016).   |
| 840        | 76. | Tamborero, D., et al. Cancer Genome Interpreter annotates the biological and clinical  |
| 841        |     | relevance of tumor alterations. <i>Genome Med</i> <b>10</b> , 25 (2018).   |
| 842        | //. | Bailey, M.H., et al. Comprehensive Characterization of Cancer Driver Genes and Mutations.  |
| 843        | 70  | <i>Cell</i> <b>173</b> , 3/1-385 e318 (2018).  |
| 844        | /8. | Tate, J.G., et al. COSMIC: the Catalogue Of Somatic Mutations in Cancer. Nucleic Acias Res 47,   |
| 845        | 70  | D941-0947 (2019).  |
| 846        | 79. | Alexandrov, L.B., et al. The repertoire of mutational signatures in numan cancer. Nature 578,  |
| 847        | 80  | 94-101 (2020).   |
| 040<br>040 | 80. | oanney, J.J., et ul. controls of nucleosome positioning in the numan genome. PLOS Genet 8,   |
| 049<br>850 | Q1  | ELUUSUSU (2012).<br>R Core Team : A Language and Environment for Statistical Computing R Equindation for   |
| 850<br>851 | 01. | Statistical Computing Vienna Austria https://www.R-project.org   |
| 001        |     | Statistical compating, Vienna, Austria, <u>Inteps//www.n-project.org</u> .   |













## ΔclbR E.coli infection



pks<sup>+</sup> E.coli infection



White: yH2AX, blue: nuclei, red: phalloidin

















Α



### Mutation enrichment for patterns





### Proportion of total mutations overlapping patterns







# **COSMIC Database Unique Mutations**



Mutation class At AAATT/AAAAT motif At AAWWTT motif At pos 2/5 of AAWWTT motif



Distance to motif start