1      **Colibactin DNA damage signature indicates causative role in colorectal cancer**

2

3      Paulina J. Dziubańska-Kusibab[1†], Hilmar Berger[1†], Federica Battistini[2], Britta A. M. Bouwman[3],

4      Amina Iftekhar[1], Riku Katainen[4], Nicola Crosetto[3], Modesto Orozco[2,5], Lauri A. Aaltonen[4], and

5      Thomas F. Meyer[1,*]

6      † Shared first authors

7

8      [1] Department of Molecular Biology, Max Planck Institute for Infection Biology, 10117 Berlin,
9      Germany

10     [2] Institute for Research in Biomedicine (IRB Barcelona). The Barcelona Institute of Science and
11     Technology, 08028 Barcelona, Spain

12     [3] Science for Life Laboratory, Department of Medical Biochemistry and Biophysics, Karolinska
13     Institutet, 17165 Stockholm, Sweden

14     [4] Applied Tumor Genomics Research Program and Department of Medical and Clinical
15     Genetics, Medicum, University of Helsinki, 00014 Helsinki, Finland

16     [5] Department of Biochemistry and Biomedicine, University of Barcelona, Barcelona, Spain

17

18     * Corresponding author:

19     Prof. Thomas F. Meyer
20     Department of Molecular Biology
21     Max Planck Institute for Infection Biology
22     Charitéplatz 1
23     10117 Berlin
24     Germany
25     Email: tfm@mpiib-berlin.mpg.de
26

27 **Abstract**

28 Colibactin, a potent genotoxin of *Escherichia coli*, causes DNA double strand breaks (DSBs) in human

29 cells. We investigated if colibactin creates a particular DNA damage signature in infected cells by

30 identifying DSBs in colon cells after infection with *pks+ E.coli*. Interestingly, genomic contexts of DSBs

31 were enriched for AT-rich penta-/hexameric sequence motifs, exhibiting a particularly narrow minor

32 groove width and extremely negative electrostatic potential. This corresponded with the binding

33 characteristics of colibactin to double-stranded DNA, as elucidated by docking and molecular dynamics

34 simulations. A survey of somatic mutations at the colibactin target sites of several thousand cancer

35 genomes revealed significant enrichment of the identified motifs in colorectal cancers. Our work

36 provides direct evidence for a role of colibactin in the etiology of human cancer.

37 **One sentence summary:** We identify a mutational signature of colibactin, which is significantly

38 enriched in human colorectal cancers.

39

40

41    The mucosal epithelium is a preferred target of damage by chronic bacterial infections and associated

42    toxins. Not surprisingly, most cancers originate from this tissue. Several infectious agents have been

43    implicated in human cancers, with *Helicobacter pylori* representing the prototype of a cancer-inducing

44    bacterium. Yet, unlike for infections with tumor viruses, which deposit telltale transforming genes in

45    infected cells, for bacterial pathogens compelling evidence of a carcinogenic function is missing due to

46    the lack of specific signatures of past infections in the emerging cancer genomes. Nonetheless, a

47    broader role of bacterial pathogens in human carcinogenesis is highly suggestive.

48    In humans, several bacterial species have been attributed to a potential role in colorectal cancer (CRC),

49    including *Fusobacterium nucleatum* [1] and colibactin-producing strains of *E. coli* [2,3]. Mechanistic analyses

50    indicated distinct cancer-promoting mechanisms elicited by these bacteria, including the activation of

51    inflammatory and growth-promoting signaling pathways as well as the induction of DNA damage [4]. In

52    particular, colibactin toxin, a secondary metabolite produced by strains of the B2 phylogenetic group

53    of *E. coli*, has long been known to possess DNA damaging ability. In 2006, Nougayrède and

54    collaborators described the 54 kilobase *pks* genomic island that encodes this polyketide-peptide hybrid

55    and showed that *pks*-harboring *E. coli* induce double-strand breaks (DSBs) in host cells and activate the

56    G2-M DNA damage checkpoint pathway [5]. The recent discovery of a cyclopropane ring, characteristic

57    of DNA alkylating agents, led to the isolation of colibactin-dependent N-3 adenine adducts from host

58    DNA [6]. This observation was followed by the resolution of colibactin's mature structure as a highly

59    symmetrical molecule, containing identical cyclopropane warheads at each end, which can give rise to

60    DNA cross-links[7]. Yet, it is unclear if colibactin's mode of action generates a specific signature that is

61    retrievable in cancers from tissues potentially exposed to respective *E.coli* infections.

62    To determine a potential preference of colibactin action for specific sites in host cell DNA, we began

63    by globally defining the occurrence of DSBs upon infection of colon derived cells with *pks+ E. coli.* To

64    this end, we applied 'Breaks Labeling In Situ and Sequencing' (BLISS), which allows the detection of the

3

65    exact sites of DSBs in fixed host cells [8]. The resulting next-generation sequencing (NGS) data and

66    computational analyses revealed a highly specific DNA damage signature, involving AT-rich sequence

67    patterns associated with extreme shape characteristics, which was confirmed by in silico modelling of

68    the colibactin interaction with DNA. By using this information for a stringent search of a mutational

69    signature of colibactin in human cancer genome data, we establish a role of colibactin in the cause of

70    human colorectal cancer and possibly additional cancer types.

71

**An unbiased sequencing approach to detect colibactin-induced DSB patterns**

73    To confirm colibactin-induced damage, we infected the human colorectal adenocarcinoma cell line

74    Caco-2 with *pks+ E. coli* at MOI 20 for 3 hours. Fluorescence immunohistochemistry showed that cells

75    infected with the wild-type bacteria (*pks+*) were positive for the DNA damage marker γH2AX, while

76    cells infected with the *clb*R deletion mutant (*pks-*), in which colibactin synthesis is restricted [9], were

77    not (Fig. 1A). To specifically gain insight into colibactin-induced DSBs and the DNA sequences at which

78    they occur, BLISS was applied to Caco-2 cells infected with WT M1/5 (*pks+*) and mutant Δ*clb*R M1/5

79    (*pks-*) *E. coli*. Untreated cells and cells treated with the DSB-inducing chemical agent etoposide served

80    as controls (Fig. 1B). BLISS enables unbiased identification of host cell DSBs on a genome-wide scale at

81    nucleotide resolution, based on the amplification of tagged DSBs by *in vitro* transcription. After

82    infection, cells were fixed and the preserved DSBs were blunted *in situ* to allow ligation of specific

83    double-stranded adapters containing a barcode, a unique molecular identifier (UMI), an RA5 Illumina

84    sequencing adapter and a T7 promoter sequence (Fig. 1B). After in vitro transcription, NGS libraries

85    were generated from the produced RNA and sequenced in single-end mode. The included UMIs are

86    used for PCR duplicate removal, while the sample barcodes allow for pooling of different samples prior

87    to the transcription reaction. The raw reads served to determine the genomic positions of DSBs as well

88    as the counts of unique cleavage events using an established analysis pipeline (see Methods). To

89    confirm that our method captured known DSB patterns, we examined the breakpoint density around

4

90   transcription start sites (TSS), which are reportedly susceptible to breaks induced by etoposide [8,10,11].

91   An increase in breakpoint counts around TSSs in our etoposide control was indeed observed (Fig. 1C),

92   indicating the reliability of BLISS as an approach to define colibactin-induced DSB patterns. Next, we

93   performed Locus Overlap Analysis (LOLA) to determine whether the identified DSBs were enriched in

94   particular genomic regions [12]. Interestingly, unlike the DSBs induced by treatment with etoposide or

95   the DSBs observed in the negative controls, those induced in the *pks+ E. coli* condition did not show

96   strong correlation with any known particular genomic regions (Fig. 1D).

97

98   **Colibactin damages DNA preferentially in specific AT-rich motifs**

99   Next, we asked whether we could identify any particular sequence pattern around the identified DSBs.

100   We thus analyzed nucleotide sequence content of different length stretches around all identified DSBs

101   and compared them between the different treatments. We found that DSBs in cells exposed to *pks+*

102   *E. coli* are enriched in AT-rich regions. This enrichment was particularly high for the pentanucleotides

103   AAATT and AAAAT together with their complementary mates (Fig. 2A, left panel). This sequence

104   preference of colibactin was evident when compared with either *pks- E. coli* infected or non-treated

105   cells used as the control samples (Fig. S1A). It was detected independently in all four biological

106   replicates, with almost identical relative enrichments (Fig. 2B, Tab. S1). Importantly, no meaningful

107   sequence enrichments were detected when sequence content in close proximity to the DSBs observed

108   in cells exposed to *pks- E. coli* was compared to that in non-treated cells (Fig. 2A, right panel). Hence,

109   the preference for AT-rich sequences is directly linked to the action of colibactin, rather than *E. coli*

110   infection *per se*. To identify the full motif, we analyzed the independent impact of 3' and 5' flanking

111   sequences in both identified pentanucleotides for strength of enrichment. Motifs with up to one

112   additional 3' adenine and/or 5' thymidine bases were enriched among breakpoints while no impact

113   was observed for more distal nucleotides (Fig. 2C). We also used discriminative motif discovery

114   (DREME)[13] between breakpoint contexts from *pks+* and *pks- E.coli* infected cells to further narrow

5

115    down the motif. The top-scoring motif was identified as AAWWTT (Fig. 2D), which contains the

116    enriched pentanucleotide patterns and is compatible with the 3'/5' extensions represented in Fig. 2C.

117    This symmetric motif indicates a requirement for distant adenines on opposing strands of the double

118    helix while the preference for central A/T nucleotides might derive from dependency on additional

119    conformational conditions.

120

121    **Preferred sites of colibactin action exhibit distinct DNA shape characteristics**

122    Small molecule DNA ligands bind preferentially through intercalation and/or contacts with the double

123    helix major or minor groove, where binding specificity is usually defined by nucleotide sequence-

124    dependent DNA shape characteristics (reviewed by Tse et al.[14]). To investigate whether colibactin has

125    specific DNA shape preferences, we carried out predictions of the shape features in the proximity of

126    each detected DSB. Remarkably, in close proximity (±8 bp) of the detected breakpoints, minor groove

127    width (MGW) exhibited reproducible deviations from the line of averaged values at positions located

128    further away from the DSBs. This was true not only for samples exposed to *pks+ E. coli*, but also for all

129    other samples (Fig. 3A). In addition, all other computationally predicted DNA shape features (helical

130    twist, propeller twist, roll and electrostatic potential) also showed deviations within 8 bp of the DSBs

131    in all samples (Fig. S2). To ensure that the specific landscape of DNA shape at the DSB position is not

132    an artefact of our data analysis approach, the same analysis was performed on 10,000 sequences

133    randomly chosen from the genome (Fig. 3A, inset). Prominent fluctuations of structural properties

134    along the sequences were only observed in close proximity to the identified breaks. Regarding the

135    specific properties of colibactin, we noticed that for all DNA shape patterns the average value at the

136    exact breakpoints in *pks+ E. coli*-infected cells was markedly different from that of all other samples

137    (Fig. 3A enlargement and S2).

138    Averaged profiles describe the superposition of potentially many underlying shape motifs, many of

139    which might be attributable to DSBs generated by processes other than colibactin and are shared

6

140   across conditions. To further explore the differences between the DNA shape parameters in the

141   individual breakpoint positions of *pks+* and *pks- E. coli* infections in an unbiased manner, we applied k-

142   means clustering as unsupervised machine learning algorithm. Assigning every set of predicted values

143   of the DNA shape characteristics for each DSB to the closest centroid of 1 out of 9 clusters

144   independently for both *pks+* and *pks- E. coli* induced DSBs, resulted in specific and unique shape

145   patterns for each cluster (Fig. S3A,B). Interestingly, a quarter of all breakpoints from both infection

146   models were assigned to the respective cluster 1 (Fig. 3B), whose profile amplitudes and pattern

147   correspond to the global profile of MGW. To gain a better overview of the sequence content of each

148   cluster, the probability for the presence of each nucleotide was computed for each position (top row

149   for each cluster). As expected, MGW dips were associated with high AT-rich content in all clusters.

150   Most of those dips correlate to short AT stretches most likely caused by periodic 10 bp spaced WW

151   dinucleotide motifs in the genome sequence associated with nucleosome positioning[15], which occur at

152   different positions in the breakpoint context and are therefore distributed to separate clusters.

153   To determine whether any cluster was unique to a particular treatment, we compared clusters for both

154   infection conditions. Indeed, cluster 9 of the *pks+ E. coli* dataset was not paired with any *pks-* cluster

155   for any of the parameters examined and showed the strongest deviations of shape parameters

156   centered at the estimated DSB position. This was true regardless of whether clusters were compared

157   separately for each predicted DNA shape parameter (Fig. S3C) or for all parameters together (Fig. 3C).

158   This confirms that the sequences in proximity to the DSBs assigned to cluster 9 of *pks+ E. coli* represent

159   a group of breaks unique to this condition. Differences between MGW means for cluster 1 and 9 show

160   how strongly AT-rich sequences influence local DNA shape (Fig. 3B).

161

162   **Colibactin's binding motif corresponds to extreme DNA shape parameter values**

163   In order to unveil the features of the DNA molecules preferred by colibactin, we analyzed the structural

164   properties of the DNA stretches close to the identified DSBs. We correlated the predicted DNA shape

7

165  parameters for the central 1-2 bps of all possible pentanucleotides (1024) with the log2-ratios of

166  pentanucleotide sequence enrichment in DSB positions caused by colibactin. Remarkably, colibactin's

167  preferred pentanucleotide sequences, d(AAATT)·d(AATTT) and d(AAAAT)·d(ATTTT), were associated

168  with the narrowest minor groove widths, with values below 3 and 3.7 Å respectively, as well as with

169  some of the most negative values for propeller twist of the central base pair and extremely negative

170  electrostatic potential (Fig. 4A). Closer inspection of the inter-base pair parameter roll revealed that

171  AAATT, the most frequent pentamer that surrounds the break point, also shows very peculiar

172  conformational characteristics (Figure 4A). The DNA stretch composed of A-tract followed by T-tract

173  tract shows that the progressive narrowing of the DNA minor groove going from the 5′ to the 3′ end is

174  correlated with low roll values. Values for the DNA stiffness descriptor (40) (k_tot, see Methods

175  section), revealed that these tracts possess high intrinsic rigidity (Fig. 4A, Table S1), making them

176  difficult to distort.

177  To obtain a more complete picture of the combined effect of the DNA shape characteristics, we

178  extended this analysis to the central 5 bp of all possible 9 bp sequences and explored the multivariate

179  space defined by all DNA shape parameters and at all positions by principal component analysis. Again,

180  enriched motifs in pks+ E. coli infected cells compared to pks- E. coli stood out as an extreme group

181  among all analyzed sequences (Fig. 4B). The data suggest that colibactin's binding preference for DNA

182  stretches with the central pentanucleotides AAATT/AAAAT is driven not only by nucleotide content

183  but also by particularly extreme values of sequence-associated DNA shape attributes like MGW and

184  electrostatic potential. To probe this, we also calculated the molecular interaction potential (MIP)

185  using $Na^+$ as probe for the most and the least preferred DNA central pentamers for colibactin binding

186  (AATTT and CTTTG respectively). The isosurfaces for the two DNA sequences (Figure 4C, blue)

187  confirmed strongly different electrostatic potential correlated with different minor groove

188  conformations, which is likely to be related to the difference in colibactin binding affinity. All these

189  observations suggest that the unusually narrow minor groove together with an inherent rigidity and a

8

190    marked electrostatic potential facilitate recognition and binding of colibactin, probably maximizing its

191    interactions with the DNA.

192    In order to explore the binding between the DNA and colibactin we built a molecular model of

193    colibactin (see Methods for details) using quantum mechanics (QM) calculations as first structural

194    guess. The optimized structures were then hydrated and subjected to molecular dynamics (MD)

195    simulations (details on parametrization are discussed in Methods) using state-of-the-art simulation

196    conditions (see Methods). Colibactin appears as a rather flexible molecule, with an average end-to-

197    end distance around 13 Å (Fig. S4). This suggests it can bind 4-5 base pairs if located along the minor

198    groove, which is supported by its structure, its preference for AT-rich sequences, and its ability to

199    attack N3. HADDOCK software [16] was used as docking engine, to obtain a putative binding mode. The

200    default scoring function was supplemented by restraints forcing the orientation of the reactive

201    cyclopropane moiety towards the N3 of the adenine. The best docking poses were manually curated

202    and subjected to MD simulations (see Methods). The final putative model shows a very stable binding

203    of colibactin to the minor groove (Figure 4D), with excellent van der Waals contacts with all the walls

204    of the groove and the cyclopropane rings pointing towards the adenines on opposite strands (Fig. 4E).

205    From the equilibrium trajectory we determined that the number of base pairs involved in the binding

206    could fluctuate between 4 and 5, depending on the orientation of the cyclopropane, and the carbon

207    alkylating the N3 of the adenines (Figure 4E, enlargement). In all cases colibactin fits perfectly into the

208    narrow minor groove of the targeted sequences and adopts a spatial arrangement that would facility

209    alkylation at N3.

210

211    **Somatic mutations at colibactin target sequences indicate role in cancerogenesis**

212    Having identified a specific nucleotide sequence associated with colibactin-induced DSBs, we

213    wondered if we could identify a specific mutational signature associated with this sequence in cancers

214    that have been experimentally connected to *pks+ E. coli* infection [2,17]. Using whole-exome sequencing

9

215   (WXS) data from colorectal cancer samples [18] (n=619) and across several cancer entities in the TCGA

216   project (https://www.cancer.gov/tcga, see Methods, n=553 colorectal cancers among 10,224 tumor

217   cases in 24 cancer types), we tested whether somatic mutations are specifically enriched at the

218   identified pentanucleotide sequences. We determined the hexanucleotide-specific mutation rate for

219   all possible hexanucleotides adjusted for their frequency in exonic regions. Given colibactin's

220   demonstrated preference for alkylation of adenines, we assessed the mutation rate for single

221   nucleotide variants (SNV) at reference bases A or T. We hypothesized that preferential binding of

222   colibactin to AAWWTT motifs (i.e. AAATTT or AATTTT/AAAATT) should increase the mutation rate at

223   these motifs compared to all other hexanucleotides with the same length and nucleotide content (i.e.

224   all remaining WWWWWW motifs). Since we observed that mutation rates at AAWWTT motifs were

225   particularly high in hypermutator samples harbouring polymerase epsilon (POLE) mutations, we

226   assessed mutation rates in cohorts defined by total SNV numbers per samples and POLE-mutated

227   samples separately. We found that mutation rates in AAWWTT motifs were enriched compared to all

228   other WWWWWW motifs in colorectal cancers in both data sets analyzed (Fig. 5A). In the TCGA pan-

229   cancer data set we also found enrichment at AAWWTT motifs in stomach cancer, uterine corpus

230   endometroid cancer and breast cancer. No enrichment was found, e.g. in head and neck squamous

231   cancer, lung adenocarcinoma and lung squamous carcinoma, while enrichment only for POLE mutated

232   cases was found in bladder cancer and cervical squamous cancer (Fig. 5B).

233   We validated the findings from WXS data in a cohort of colorectal cancer assessed by whole genome

234   sequencing (WGS) [19]. We analyzed enrichment of mutations at colibactin associated motifs for 208

235   tumors including 193 microsatellite stable (MSS), 3 POLE mutated and 12 microsatellite instable (MSI)

236   cases in a similar way as for WXS data but considering each sample separately instead of pooling in

237   subcohorts. This allowed to identify enrichment and mutational loads for individual samples. We found

238   significant (Mann-Whitney-U test, p<0.05, FDR <20%) enrichment of mutations at colibactin associated

239   pentanucleotide motifs compared to other motifs with same length and A/T content in 3/3 POLE

240   mutated samples and 49/193 (25.3%) MSS cases but not in MSI cases. We found similar enrichment as

10

241  for penta- (AAATT/AAAAT) for hexanucleotide (AAWWTT) motifs associated with colibactin in MSS

242  samples (data not shown). The median number of mutations in MSS samples at colibactin associated

243  motifs was 963 (range: 63-11876) corresponding to a median proportion of 6.7% (range: 3.9-44.7%).

244  We next asked if an association exists between the preferred colibactin motif and any of the previously

245  described mutational signatures[20,21]. Again, we used somatic mutation data from the TCGA data set as

246  above and classified all single nucleotide variants according to the sequence context in direct proximity

247  (+/- 5bp). Variants were assigned to one of three groups: Those with sequence context containing

248  AAATT/AATTT or AAAAT/ATTTT, those with contexts containing a control TTT motif and all remaining

249  mutations. Globally we observed distinct mutation frequencies for several trinucleotide changes (Fig.

250  5C) in those classes. We identified a contribution of known signatures in those 3 classes for all samples

251  and selected those with significantly higher contributions at AATTT or AAAT motifs compared to TTT

252  and all other motifs (Fig. 5D). Two of the signatures with increased contributions at colibactin-

253  associated           motifs           where           of           particular           interest:           SBS28

254  (https://cancer.sanger.ac.uk/cosmic/signatures/SBS/SBS28.tt)                    and                    SBS41

255  (https://cancer.sanger.ac.uk/cosmic/signatures/SBS/SBS41.tt) both with unknown etiology and

256  featuring predominantly mutations at T:A and a prominent T[T>G]T trinucleotide change, were found

257  enriched in colorectal cancers, among others. While SBS28 has been previously shown to be associated

258  with POLE mutation-related hypermutated tumors, SBS41 was enriched in stomach adenocarcinoma,

259  colorectal adenocarcinoma and endometrial carcinoma of the uterine corpus, mirroring the results for

260  motif enrichment above.

261

262  **Discussion**

263  We pursued an unbiased bimodal approach that revealed a signature of the bacterial genotoxin

264  colibactin in the human cancer genome indicating a causal link between a bacterial infection and the

265  emergence of cancer. This was achieved by first defining the DSB-landscape generated by the action

11

266   of colibactin through applying the BLISS sequencing technology and subsequent comprehensive

267   analysis of the genome-wide location of DSBs. The resulting DSB pattern, which exhibits exceptional

268   structural features, corresponded to, and could be further refined by, three-dimensional modelling of

269   the colibactin–DNA complex, involving distinct topological interactions with the minor-groove. In a

270   second step, we used the identified motif to assign associated mutations in various cancer genome

271   databases. Most interestingly, we revealed an enrichment of mutations at colibactin-associated motifs

272   in colorectal cancers but also detectable in a few other cancer types, notably uterine endometroid and

273   stomach cancer. We identified putative trinucleotide signatures (SBS41, SBS28) in the context of these

274   mutant sites in the same cancer entities.

275   The identified AAATT and AAAAT motifs are associated with extreme physical values of the DNA

276   duplex, most prominently characterized by a very narrow minor groove width, which generates highly

277   negative electrostatic potential and renders the DNA segment stiff. This extreme physical property

278   implicates a low propensity of the colibactin target site to bind to proteins [22]. In fact, poly(dA-dT)-tracts

279   are rarely found inside nucleosomes, but are prevalent in nucleosome-free regions (NFRs) [22]. Thus,

280   colibactin's particular targeting preferences for non-protected DNA regions might increase the efficacy

281   of the toxin. Even though definitive evidence of the binding conformation requires further

282   experimental support from 3D structural analysis, the 3D model provided here allows for

283   demonstrating and validating the extraordinary electrostatic properties of the identified motifs and

284   the fit of colibactin to the minor groove. It puts a limit of 4-5 nucleotides on the distance of adenines

285   attacked by the cyclopropane groups of the same molecule. Although most of the DNA shape

286   characteristics are directly driven by the underlying sequence, the fact that other sequences with

287   similar A/T content were not strongly enriched around DSBs indicates a dominance of DNA shape over

288   sequence characteristics for the binding of colibactin.

289   Similar DNA shape and sequence affinities have been reported for other bacterial DNA toxins, such as

290   duocarmycin, yatakemycin, distamycin, netropsin and CC-1065 – small molecules produced by

291   *Streptomyces* spp., which are all minor groove binders with AT-rich sequence selectivity. Distamycin [23]

12

292   and netropsin [24] act as RNA and DNA polymerase inhibitors [25], while duocarmycin, yatakemycin and

293   CC-1065 are DNA alkylators. Duocarmycin [26] selectively alkylates adenine residues flanked by three 5'-

294   A or T-bases (5'-WWWA-3') [27], yatakemycin [28] preferentially alkylates the central adenine of a five-base

295   AT site (5'-WWAWW-3') [29] and CC-1065 [30,31] shows selectivity for more extended five-base AT-rich

296   alkylation sites (5'-WWWWA-3') [27]. The fact that these toxins possess similar mechanisms of action,

297   even though they derive from different bacterial strains, suggests that they arose via convergent

298   evolution. Genotoxins are widespread amongst bacterial species, where they are thought to serve

299   primarily for inter-microbial competition [32]. Unsurprisingly, therefore, all of the mentioned alkylating

300   toxins inhibit the growth of many Gram-positive and Gram-negative bacteria as well as some

301   pathogenic fungi, such as *Aspergillus fumigatus* and *Candida albicans* [28,30]. Similarly, *pks+ E.coli* inhibit

302   the growth of *Staphylococcus aureus*, also in its multi-resistant form [33].

303   How colibactin-induced DNA damage is repaired is still unknown. Different host DNA repair

304   mechanisms can be involved depending possibly also on the cell cycle phase. Effects of repair involve

305   nucleotide excision [34] of alkylated adenines which could lead to DSBs, resection of break ends or

306   complete repair, or error-prone repair by translesion DNA polymerases in late phases of the cell cycle,

307   among others. We were able to show enrichment of SNV at colibactin-associated motifs in exome and

308   whole-genome sequencing datasets. For colorectal cancers, whole-genome sequences revealed

309   elevated mutation rates in colibactin associated motifs in at least 25% of all MSS cases and a colibactin

310   attributable mutation load of around 6% in most patients. Further analyses of whole-genome

311   sequenced samples including the analysis of breakpoints of structural variants will be required to

312   assess the full spectrum of damage-related mutations in host cells. Mutational signatures for other

313   alkylating substances, such as cisplatin, have been identified in human DNA sequences after exposure

314   to the mutagen [21,35]. However, it is to be expected that the signatures depend strongly on the specific

315   type of damage induced by each substance. Here we identified two signatures that are consistent with

316   colibactin action, one with (SBS28) and one without (SBS41) relation to known DNA repair defects. An

317   impact of reduced DNA repair and mutagen-induced damage on the emergence of different

13

318    mutational signatures has recently been shown in a model of *C. elegans* [36]. The enrichment of

319    mutations specifically in POLE cases hints at either a similar outcome of distinct mutational processes

320    or even a role of POLE in the repair of colibactin-associated damage.

321    Colibactin has been found not only in *E. coli* but also in *Klebsiella* isolates [37]. Considering the widespread

322    and diversity of bacteria carrying this toxin, it is maybe not surprising that the mutational signature

323    identified here is not only restricted to the colon. Rather, other tissues might also be colonized by

324    either *pks*+ E. coli, another species bearing the *pks* gene cluster, or a different species with a closely

325    related genotoxin. Thus, our study will stimulate future research on other pathogen-host cell

326    encounters that could lead to an even greater match of the identified signature with different cancer

327    types. Better understanding of the role of the microbiome in malignant degeneration should provide

328    new and exciting opportunities for cancer prevention.

329

330

331    **Materials and Methods**

332    **Cell line, bacterial strains, E.coli infection and etoposide-treatment**

333    Caco-2 cells (from ATCC® HTB-37$^{TM}$) were cultured at 37 °C under a water-saturated 5% $CO_2$

334    atmosphere, in DMEM medium (Life Technologies, cat. number: 10938-025), supplemented with 20%

335    FCS (Biochrom, cat. number: S0115). Contamination of Mycoplasma spp. in immortal cell line was

336    excluded using Venor®GeM OneStep PCR kit (Minerva Biolabs®, cat. number: 11-8250). To infect Caco-

337    2 cells, overnight liquid culture of E.coli strain M1/5 (Streptomycin-resistant and colibactin-positive)

338    and E.coli strain M1/5::ΔclbR (streptomycin-resistant and colibactin-negative) was set up. Bacteria

339    were inoculated in 5 ml of Luria broth (LB) medium and incubated overnight at 37 °C in a shaking

340    incubator. The overnight inoculum was diluted 1:33 in infection medium (DMEM + 10% FCS + HEPES

341    (Life Technologies, cat. number: 15630-056)) to obtain $OD_{600}=1$ after 3 h of incubation to give $1.5 \times 10^9$

14

342     bacteria/ml. Prepared bacteria inoculum was further diluted to reach MOI 20, added to Caco-2 cells

343     seeded previously and incubated for 3 hours at 37 °C. Medium was then aspirated and cells fixed

344     according to the protocol for immunofluorescence or BLISS. For every biological replicate positive

345     (etoposide-treatment) and negative (no treatment) controls were included. Etoposide powder (Sigma

346     Aldrich, cat. number: E1383) was diluted in DMSO in order to reach 50 mM working solution. Aliquots

347     of the drug were stored at -20 °C. Final drug dilutions to the concentration of 50 µM were performed

348     in pre-warmed infection medium prior to each drug exposure. Treatment was conducted for 3 hours

349     at 37 °C and afterwards medium was aspirated and etoposide-treated cells were fixed in the same way

350     as E. coli-infected cells.

351     **Immunofluorescence staining**

352     Caco-2 cells grown and infected on MatTek glass-bottom dishes were washed three times with PBS

353     (Life Technologies, cat. number: 14190-094) and fixed with 3.7% paraformaldehyde (Sigma Aldrich,

354     cat. number: P6148) for 1 h. The cells were kept overnight in blocking buffer (3% BSA, Biomol, cat.

355     number: 01400.100), 1% saponin (Sigma Aldrich, cat. number: 84510), 2% Triton X-100 (Carl Roth, cat.

356     number: 3051.2) and 0.02% sodium azide (Sigma Aldrich, cat. number: S2002). Blocking was followed

357     by overnight incubation with γH2AX antibody (Phospho-Histone H2A.X (Ser139) Antibody, Cell

358     Signaling, cat. number:  2577, 1:500 dilution) at 4 °C. The next day, the MatTek dishes were washed

359     three times with blocking buffer followed by overnight incubation with secondary antibody (Dianova,

360     cat. number: 711-035-152, 1:250 dilution) diluted in blocking buffer. Phalloidin 546 (Invitrogen, cat.

361     number: A22283, 1:200 dilution) and Hoechst (Sigma, cat. number: H6024, 1:10000 dilution) were

362     added for staining actin filaments and DNA, respectively. The next day, cells were washed three times

363     with blocking buffer and coverslipped using Vectashield® Antifade Mounting Medium (Vector

364     Laboratories, cat. number: H-1000). Images were acquired using a Leica TCS SP-8 confocal microscope

365     and processed using ImageJ.

366

367     **sBLISS, an adaptation of the BLISS method**

15

368    DSBs were identified using the suspension-cell BLISS (sBLISS) method [38], which is an adaptation of the

369    previously published BLISS protocol [8,39]. In contrast to BLISS, where DSBs are labeled in fixed cells

370    immobilized on microscope slide, in sBLISS DSBs are labeled in fixed cell suspensions. In brief, cells

371    were treated/infected in culture dishes and afterwards trypsinized, counted, centrifuged and

372    resuspended in pre-warmed medium to obtain $10^6$ cells per 1 ml. Then, cells were fixed by adding 16%

373    PFA (Electron Microscopy Sciences, cat. number: 15710) to reach a final concentration of 4%. After 10

374    minutes, 2 M glycine (Molecular Dimensions, cat. number: MD2-100-105) was added to a final

375    concentration of 125 mM in order to block unreacted aldehydes. This was followed by two 5 minutes

376    incubations, first at room temperature and then on ice, followed by two washes in ice-cold PBS. Cross-

377    linked cells were stored in PBS at 4 °C until further processing.

378    Next, BLISS template was prepared. This includes: (1) Cell lysis in 10mM Tris-HCl, 10 mM NaCl, 1 mM

379    EDTA, and 0.2% Triton X-100 (pH 8) buffer, followed with lysis in buffer containing 10 mM Tris-HCl, 150

380    mM NaCl, 1 mM EDTA, and 0.3% SDS (pH 8); (2) DSBs blunting with NEB's Quick Blunting Kit (NEB, cat.

381    number: E1201); (3) *In situ* BLISS adapter ligation using T4 DNA Ligase (ThermoFisher Scientific, cat.

382    number: EL0011). Each BLISS adapter contained a T7 promoter sequence for IVT, the RA5 Illumina RNA

383    adapter sequence, a random 8nt long sequence referred to as Unique Molecular Identifier (UMI) and

384    a 8nt long sample barcode; (4) Phenol:chloroform-based extraction of gDNA; (5) Fragmentation of

385    isolated genomic DNA (400-600bp) using BioRuptor Plus (Diagenode). Obtained BLISS templates were

386    stored at -20 °C.

387    The final step of the BLISS protocol was *in vitro* transcription (IVT) followed by NGS library preparation.

388    At first, 100ng of purified, sonicated and differentially-barcoded BLISS template of 1) etoposide-

389    treated and non-treated cells, or 2) cells infected with pks+ E.coli or infected with pks- E.coli were

390    pooled into one reaction, respectively. IVT was performed using MEGAscript T7 Transcription Kit

391    (ThermoFisher, cat. number: AMB13345) for 14 hours at 37 °C in the presence of RiboSafe RNAse

392    Inhibitor (Bioline, cat. number BIO-65028). Next, gDNA was removed using DNase I (ThermoFisher, cat.

393    number: AM2222) and the remaining RNA was purified with Agencourt RNAClean XP beads (Beckman

16

394    Coulter). The Illumina RA3 adapter sequence was ligated to the purified RNA using T4 RNA Ligase 2

395    (NEB, cat. number: M0242) for 2 hours at 25 °C and reverse transcription was performed with Reverse

396    Transcription Primer (Illumina sequence) using SuperScript IV Reverse Transcriptase (ThermoFisher,

397    cat. number: 18090050) for 50 minutes at 50 °C. This was followed by enzyme heat inactivation for 10

398    minutes at 80 °C. Finally, libraries were amplified with NEBNext High-Fidelty 2x PCR Master Mix (NEB,

399    cat. number: M0541), the RP1 common primer and a uniquely selected index primer. 12 PCR cycles

400    were conducted, and after that libraries were purified according to the two-sided AMPure XP bead

401    purification protocol (Beckman Coulter). Profiles of the libraries were quantified on a BioAnalyzer High

402    Sensitivity DNA chip. Libraries were sequenced as single-end (1x75) reads on the NextSeq platform.

403    **Pre-processing of sequencing data**

404    Raw sequencing data were pre-processed as previously described [7]. In brief, only reads which

405    contained the expected prefix of UMI and sample barcode were kept using SAMtools [40]. One mismatch

406    in the barcode sequence was allowed. Further, prefixes were trimmed and the remaining sequences

407    were aligned to the GRCh37/hg19 reference genome using BWA-MEM [41]. Reads with mapping quality

408    scores ≤ 30 and those which were determined as PCR duplicates were removed. Finally, a BED file

409    containing a list of unique DSBs locations was generated. DSBs which fell into ENCODE blacklist regions

410    [42], high coverage regions [34] and low mappability regions [34] were removed. Kept positions of DSBs were

411    further used in downstream analysis.

412    **Locus Overlap Analysis**

413    To identify significant overlaps of DNA DSB with genomic region sets we used LOLA [11]. We first defined

414    whole genome as a Universe Set, which was next divided into tiles of equal lengths (1,000 nt). For each

415    created tile we next searched for overlaps with captured by BLISS DSBs using the findOverlap()

416    function. All tiles containing ≥ 10 breaks were used as a Query Set. The runLOLA() function was

417    executed with LOLA Core databases (reduced by Tissue clustered DNase hypersensitive sites) as well

17

418    as LOLA Extended databases and custom database containing non-B-DNA regions (https://nonb-

419    abcc.ncifcrf.gov/apps/site/references). Fisher's exact test was used with a FDR ≤ 5%.

**DNA Shape predictions**

421    DNA structures can be described in terms of base-pair and base-step parameters that consist of three

422    translational and rotational movements between the bases or the base pairs, respectively. At the base-

423    pair step level, DNA deformability along these six directions has been described by the associated

424    stiffness matrix [43]. From the ensemble of MD simulations considering the tetramer environment using

425    the newly refined parmbsc1 force field, we retrieved the 6x6 matrix describing the deformability of

426    the helical parameters for each possible DNA tetramer. Pure stiffness constants corresponding to the

427    six base-pair step parameters (shift, slide, rise, tilt, roll and twist) were extracted from the diagonal of

428    the matrix and the total stiffness ($K\_$tot) was obtained as a product of these six constants and used as

429    an estimate of the flexibility of each base pair step in a tetramer. For predictions of minor groove width

430    (MGW), propeller twist (ProT), electrostatic potential (EP), helical twist (HelT) and roll (Roll) the

431    getShape function from 'DNAshapeR' package was used [44]. Input FASTA files, containing sequences in

432    close proximity to identified DSB (±5nt or ±100nt), were extracted with custom python script (available

433    upon request). The interaction potential (electrostatic and van der Waals) of $Na^+$ probes with DNA

434    duplexes was determined using a linear approximation to the Poisson-Boltzmann equation and

435    dielectric constant for the DNA as implemented in the CMIP program [45].

**K-means clustering of DNA shape profiles**

437    We used an elbow method to find appropriate number of clusters in the dataset, which consisted of

438    predicted values of all parameters (MGW, HelT, ProT, Roll, EP) ±8 nt from each breakpoint. Based on

439    cluster number diagnostic it was chosen to use k=9. Initial cluster centers were defined using 100

440    iterations. Next, we assigned every set of observations for each breakpoint into the closest centroid of

441    1 out of 9 clusters, independently for both – pks+ and pks- E.coli-induced DSB. Finally, sequence

442    content of each cluster was exported and used as an input for computing proportion of each nucleotide

443    per position (see SeqLogo method).

444    **SeqLogo**

445    To compute and visualize the proportion of each nucleotide per position from collection of sequences

446    consensusMatrix() and seqLogo() functions from 'seqLogo' package were used [46,47]

447    **Model and Molecular dynamics set up**

448    The 3D structure and protonation state of the colibactin were built starting from the smile

449    (https://pubchem.ncbi.nlm.nih.gov/compound/138805674#section=InChI)    using    MarvinSketch

450    (MarvinSketch,    version    6.2.2,    calculation    module    developed    by    ChemAxon,

451    http://www.chemaxon.com/products/marvin/marvinsketch/). The geometry of the model and the

452    partial atomic charges were assigned to the structure with General Amber Force Field (GAFF) [48].

453    Parameters and topology files were prepared with Acpype [49]. The colibactin was then simulated in

454    explicit solvent at 298K (see below for details) for 250ns and along the simulation the distance between

455    the cyclopropanes was monitored (see Fig.S4), to study their orientation and the overall length of the

456    free colibactin. Using HADDOCK 2.4 [16], we then built the complex DNA-colibactin. For the docking, we

457    selected a representative structure of the free colibactin along the MD simulation, with an average

458    distance among the cyclopropanes (red line, Fig. S4) and an equilibrated structure of the DNA

459    (sequence CGAAATTTCG). After the initial docking, that positioned the molecule correctly along the

460    minor groove of the DNA, we then manually rotated slightly the molecule to improve the orientation

461    of the cyclopropanes towards the N3 of the closest adenine using PYMOL (The PyMOL Molecular

462    Graphics System, Schrödinger, LLC (2018)). To check the stability of this complex and to equilibrate its

463    structure the model was simulated (see details MD simulation below) and minimized in solution with

464    positional restraints on the solute using our well-established multi-step protocol [50,51]. The minimized

465    structure was thermalized to 298K at NVT, and then simulated first applying harmonic restraints of 5

466    kcal/mol·Å2 on the DNA on the DNA structure and distance constraints between the cyclopropane and

19

467     the N3 of the adenine (respectively 4 and 5 bases apart), each represented by a harmonic restraint of

468     2.5 kcal/mol·$\text{Å}^2$. To further check the stability of the complex we then slowly removed the constraints

469     and run MD simulation of the complex during 60 ns by means of Molecular Dynamics simulations at

470     NPT (P = 1 atm; T= 298K). The first 10 ns of the simulations were considered as an equilibration step

471     and were discarded for further analysis.

472     In each MD simulation, DNA, free colibactin and their complex, respectively, we placed the solute in

473     the centre of a truncated octahedral box of TIP3P water molecules [52], neutralized by K+ ions. In each

474     simulation the Berendsen algorithm [53] was used to control the temperature and the pressure, with a

475     coupling constant of 5 ps; and the SHAKE algorithm was utilized to equilibrium the length of hydrogen

476     atoms involved in the covalent bonds [54]. Long-range electrostatic interactions were accounted for by

477     using the Particle Mesh Ewald method (14) with standard defaults, and a real-space cut-off of 10 Å.

478     For the DNA we used the newly revised force field parmBSC1 [55]. All simulations were carried out using

479     AMBER 18 [56], and analyzed with CPPTRAJ [57] and visualized using VMD 1.9.4 [58].

480

481     **Cancer somatic mutation data**

482     We obtained somatic variant data from the TCGA Unified Ensemble "MC3" Call Set [59] ("TCGA pan-

483     cancer dataset") and from the supplementary data of Giannakis et al [18]. To test for enrichment of

484     mutations at any motif we first identified positions of all hexanucleotide motifs in the exonic portion

485     of the genome. Somatic variants occurring at A or T bases were grouped in one of 6 classes (quartile

486     1-4, outlier or POLE mutated sample) depending on the total SNV number and POLE mutation status

487     of the corresponding tumor sample We then computed the mutation rate for each hexanucleotide

488     motif with respect to the number of genomic bases covered in exonic regions for the same motif. As a

489     baseline, we established the mutation rates of all WWWWWW motifs and subtracted their mean from

490     the mutation rate of all other hexanucleotide motifs. We then tested for significance of the mutation

491     rate at colibactin associated AAWWTT motifs (i.e. AAATTT and AAAATT/AATTTT) compared to the

20

492     remaining WWWWWW motifs using Mann-Whitney-U tests and computed the false discovery rate

493     (FDR) using the method of Benjamini-Hochberg [60]. Reads from WGS of colorectal cancers [19] EGA

494     database accession code EGAS00001003010, ) were aligned to GRCh38 with BWA-MEM [41] and called

495     using Mutect2 [61]. All single nucleotide variant calls (PASSed by Mutect2) were used to determine the

496     number of mutations overlapping WWWWW pentanucleotides and WWWWWW hexanucleotides and

497     further analyzed in a similar way as for exome sequencing data on an individual sample basis.

498     **Analysis of pattern enrichment in cancers**

499     For analysis of signatures we classified all variants according to the presence of patterns in the +/- 5bp

500     around SNV variant calls: one group contained colibactin associated pentanucleotides (AAATT/AATTT

501     or AAAAT/ATTTT), one contained AAA/TTT in order to control for AT-rich sequences and one contained

502     all other motifs. The R package deconstructSigs [62] was used to estimate the contribution of COSMIC

503     signatures v3 [21] independently for each group. Differences between groups were assessed for each

504     single base change signature (SBS) between groups using Mann-Whitney test.

505     **Data analysis and visualization**

506     All visualizations and statistical analyses were produced using R v3.4 [63]

507

508

509 **References and Notes**

510 1    Castellarin, M. *et al.* Fusobacterium nucleatum infection is prevalent in human colorectal
511      carcinoma. *Genome Res* **22**, 299-306, doi:10.1101/gr.126516.111 (2012).
512 2    Arthur, J. C. *et al.* Intestinal inflammation targets cancer-inducing activity of the microbiota.
513      *Science* **338**, 120-123, doi:10.1126/science.1224820 (2012).
514 3    Cougnoux, A. *et al.* Bacterial genotoxin colibactin promotes colon tumour growth by inducing
515      a senescence-associated secretory phenotype. *Gut* **63**, 1932-1942, doi:10.1136/gutjnl-2013-
516      305257 (2014).
517 4    Bleich, R. M. & Arthur, J. C. Revealing a microbial carcinogen. *Science* **363**, 689-690,
518      doi:10.1126/science.aaw5475 (2019).
519 5    Nougayrede, J. P. *et al.* Escherichia coli induces DNA double-strand breaks in eukaryotic cells.
520      *Science* **313**, 848-851, doi:10.1126/science.1127059 (2006).
521 6    Wilson, M. R. *et al.* The human gut bacterial genotoxin colibactin alkylates DNA. *Science* **363**,
522      doi:10.1126/science.aar7785 (2019).
523 7    Xue, M. *et al.* Structure elucidation of colibactin and its DNA cross-links. *Science* **365**,
524      doi:10.1126/science.aax2685 (2019).
525 8    Yan, W. X. *et al.* BLISS is a versatile and quantitative method for genome-wide profiling of
526      DNA double-strand breaks. *Nat Commun* **8**, 15058, doi:10.1038/ncomms15058 (2017).
527 9    Brotherton, C. A., Wilson, M., Byrd, G. & Balskus, E. P. Isolation of a metabolite from the pks
528      island provides insights into colibactin biosynthesis and activity. *Org Lett* **17**, 1545-1548,
529      doi:10.1021/acs.orglett.5b00432 (2015).
530 10   Canela, A. *et al.* Genome Organization Drives Chromosome Fragility. *Cell* **170**, 507-521.e518,
531      doi:10.1016/j.cell.2017.06.034 (2017).
532 11   Yang, F., Kemp, C. J. & Henikoff, S. Anthracyclines induce double-strand DNA breaks at active
533      gene promoters. *Mutat Res* **773**, 9-15, doi:10.1016/j.mrfmmm.2015.01.007 (2015).
534 12   Sheffield, N. C. & Bock, C. LOLA: enrichment analysis for genomic region sets and regulatory
535      elements in R and Bioconductor. *Bioinformatics* **32**, 587-589,
536      doi:10.1093/bioinformatics/btv612 (2016).
537 13   Bailey, T. L. DREME: motif discovery in transcription factor ChIP-seq data. *Bioinformatics* **27**,
538      1653-1659, doi:10.1093/bioinformatics/btr261 (2011).
539 14   Tse, W. C. & Boger, D. L. Sequence-selective DNA recognition: natural products and nature's
540      lessons. *Chem Biol* **11**, 1607-1617, doi:10.1016/j.chembiol.2003.08.012 (2004).
541 15   Segal, E. *et al.* A genomic code for nucleosome positioning. *Nature* **442**, 772-778,
542      doi:10.1038/nature04979 (2006).
543 16   van Zundert, G. C. P. *et al.* The HADDOCK2.2 Web Server: User-Friendly Integrative Modeling
544      of Biomolecular Complexes. *J Mol Biol* **428**, 720-725, doi:10.1016/j.jmb.2015.09.014 (2016).
545 17   Buc, E. *et al.* High prevalence of mucosa-associated E. coli producing cyclomodulin and
546      genotoxin in colon cancer. *PLoS One* **8**, e56964, doi:10.1371/journal.pone.0056964 (2013).
547 18   Giannakis, M. *et al.* Genomic Correlates of Immune-Cell Infiltrates in Colorectal Carcinoma.
548      *Cell Rep* **15**, 857-865, doi:10.1016/j.celrep.2016.03.075 (2016).
549 19   Katainen, R. *et al.* CTCF/cohesin-binding sites are frequently mutated in cancer. *Nat Genet*
550      **47**, 818-821, doi:10.1038/ng.3335 (2015).
551 20   Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res* **47**,
552      D941-d947, doi:10.1093/nar/gky1015 (2019).
553 21   Alexandrov, L. B. *et al.* The Repertoire of Mutational Signatures in Human Cancer. *bioRxiv*,
554      322859, doi:10.1101/322859 (2018).
555 22   Chakraborty, A. *et al.* DNA structure directs positioning of the mitochondrial genome
556      packaging protein Abf2p. *Nucleic Acids Res* **45**, 951-967, doi:10.1093/nar/gkw1147 (2017).
557 23   Arcamone, F., Penco, S., Orezzi, P., Nicolella, V. & Pirelli, A. STRUCTURE AND SYNTHESIS OF
558      DISTAMYCIN A. *Nature* **203**, 1064-1065, doi:10.1038/2031064a0 (1964).

22

24  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).

25  Hahn, F. E. Distamycins and netropsin as inhibitors of RNA and DNA polymerases. *Pharmacology & Therapeutics. Part A: Chemotherapy, Toxicology and Metabolic Inhibitors* **1**, 475-485 (1977).

26  Takahashi, I. *et al.* Duocarmycin A, a new antitumor antibiotic from Streptomyces. *J Antibiot (Tokyo)* **41**, 1915-1917, doi:10.7164/antibiotics.41.1915 (1988).

27  Boger, D. L. & Johnson, D. S. CC-1065 and the duocarmycins: unraveling the keys to a new class of naturally derived DNA alkylating agents. *Proc Natl Acad Sci U S A* **92**, 3642-3649, doi:10.1073/pnas.92.9.3642 (1995).

28  Igarashi, Y. *et al.* Yatakemycin, a novel antifungal antibiotic produced by Streptomyces sp. TP-A0356. *J Antibiot (Tokyo)* **56**, 107-113, doi:10.7164/antibiotics.56.107 (2003).

29  Parrish, J. P., Kastrinsky, D. B., Wolkenberg, S. E., Igarashi, Y. & Boger, D. L. DNA alkylation properties of yatakemycin. *J Am Chem Soc* **125**, 10971-10976, doi:10.1021/ja035984h (2003).

30  Hanka, L. J., Dietz, A., Gerpheide, S. A., Kuentzel, S. L. & Martin, D. G. CC-1065 (NSC-298223), a new antitumor antibiotic. Production, in vitro biological activity, microbiological assays and taxonomy of the producing microorganism. *J Antibiot (Tokyo)* **31**, 1211-1217, doi:10.7164/antibiotics.31.1211 (1978).

31  Hurley, L. H. & Rokem, J. S. Biosynthesis of the antitumor antibiotic CC-1065 by Streptomyces zelensis. *J Antibiot (Tokyo)* **36**, 383-390, doi:10.7164/antibiotics.36.383 (1983).

32  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, doi:10.1038/nrmicro2259 (2010).

33  Fais, T. *et al.* Antibiotic Activity of Escherichia coli against Multiresistant Staphylococcus aureus. *Antimicrob Agents Chemother* **60**, 6986-6988, doi:10.1128/aac.00130-16 (2016).

34  Martin, L. P., Hamilton, T. C. & Schilder, R. J. Platinum resistance: the role of DNA repair pathways. *Clin Cancer Res* **14**, 1291-1295, doi:10.1158/1078-0432.Ccr-07-2238 (2008).

35  Boot, A. *et al.* In-depth characterization of the cisplatin mutational signature in human cell lines and in esophageal and liver tumors. *Genome Res* **28**, 654-665, doi:10.1101/gr.230219.117 (2018).

36  Volkova, N. V. *et al.* Mutational signatures are jointly shaped by DNA damage and repair. *bioRxiv*, 686295, doi:10.1101/686295 (2019).

37  Putze, J. *et al.* Genetic structure and distribution of the colibactin genomic island among members of the family Enterobacteriaceae. *Infect Immun* **77**, 4696-4703, doi:10.1128/iai.00522-09 (2009).

38  Gothe, H. J. *et al.* Spatial Chromosome Folding and Active Transcription Drive DNA Fragility and Formation of Oncogenic MLL Translocations. *Mol Cell* **75**, 267-283.e212, doi:10.1016/j.molcel.2019.05.015 (2019).

39  Zhang, F. *et al.* Breaks Labeling in situ and sequencing (BLISS). *Protocol Exchange* **DOI: 10.1038/protex.2017.018** (2017).

40  Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).

41  Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).

42  An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74, doi:10.1038/nature11247 (2012).

43  Drsata, T. *et al.* Mechanical properties of symmetric and asymmetric DNA A-tracts: implications for looping and nucleosome positioning. *Nucleic Acids Res* **42**, 7383-7394, doi:10.1093/nar/gku338 (2014).

44  Chiu, T. P. *et al.* DNAshapeR: an R/Bioconductor package for DNA shape prediction and feature encoding. *Bioinformatics* **32**, 1211-1213, doi:10.1093/bioinformatics/btv735 (2016).

45  Gelpi, J. L. *et al.* Classical molecular interaction potentials: improved setup procedure in molecular dynamics simulations of proteins. *Proteins* **45**, 428-437 (2001).

23

611  46  Schneider, T. D. & Stephens, R. M. Sequence logos: a new way to display consensus
612      sequences. *Nucleic Acids Res* **18**, 6097-6100, doi:10.1093/nar/18.20.6097 (1990).
613  47  Bembom, O. seqLogo: Sequence logos for DNA sequence alignments. R package version
614      1.44.0.  (2017).
615  48  Wang, J., Wolf, R. M., Caldwell, J. W., Kollman, P. A. & Case, D. A. Development and testing of
616      a general amber force field. *J Comput Chem* **25**, 1157-1174, doi:10.1002/jcc.20035 (2004).
617  49  Sousa da Silva, A. W. & Vranken, W. F. ACPYPE - AnteChamber PYthon Parser interfacE. *BMC*
618      *Res Notes* **5**, 367, doi:10.1186/1756-0500-5-367 (2012).
619  50  Dans, P. D. *et al.* Long-timescale dynamics of the Drew-Dickerson dodecamer. *Nucleic Acids*
620      *Res* **44**, 4052-4066, doi:10.1093/nar/gkw264 (2016).
621  51  Perez, A., Luque, F. J. & Orozco, M. Dynamics of B-DNA on the microsecond time scale. *J Am*
622      *Chem Soc* **129**, 14739-14745, doi:10.1021/ja0753546 (2007).
623  52  Jorgensen, W. L., Chandrasekhar, J., Madura, J. D., Impey, R. W. & Klein, M. L. Comparison of
624      simple potential functions for simulating liquid water. *The Journal of chemical physics* **79**,
625      926-935 (1983).
626  53  Berendsen, H. J., Postma, J. v., van Gunsteren, W. F., DiNola, A. & Haak, J. R. Molecular
627      dynamics with coupling to an external bath. *The Journal of chemical physics* **81**, 3684-3690
628      (1984).
629  54  Ryckaert, J.-P., Ciccotti, G. & Berendsen, H. J. Numerical integration of the cartesian
630      equations of motion of a system with constraints: molecular dynamics of n-alkanes. *Journal*
631      *of computational physics* **23**, 327-341 (1977).
632  55  Ivani, I. *et al.* Parmbsc1: a refined force field for DNA simulations. *Nat Methods* **13**, 55-58,
633      doi:10.1038/nmeth.3658 (2016).
634  56  Case, D. *et al. AMBER 18. University of California, San Francisco* (2018).
635  57  Roe, D. R. & Cheatham, T. E., 3rd. PTRAJ and CPPTRAJ: Software for Processing and Analysis
636      of Molecular Dynamics Trajectory Data. *J Chem Theory Comput* **9**, 3084-3095,
637      doi:10.1021/ct400341p (2013).
638  58  Humphrey, W., Dalke, A. & Schulten, K. VMD: visual molecular dynamics. *J Mol Graph* **14**, 33-
639      38, 27-38 (1996).
640  59  Ellrott, K. *et al.* Scalable Open Science Approach for Mutation Calling of Tumor Exomes Using
641      Multiple Genomic Pipelines. *Cell Syst* **6**, 271-281.e277, doi:10.1016/j.cels.2018.03.002 (2018).
642  60  Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful
643      approach to multiple testing. *Journal of the Royal statistical society: series B*
644      *(Methodological)* **57**, 289-300 (1995).
645  61  Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and
646      heterogeneous cancer samples. *Nat Biotechnol* **31**, 213-219, doi:10.1038/nbt.2514 (2013).
647  62  Rosenthal, R., McGranahan, N., Herrero, J., Taylor, B. S. & Swanton, C. DeconstructSigs:
648      delineating mutational processes in single tumors distinguishes DNA repair deficiencies and
649      patterns of carcinoma evolution. *Genome Biol* **17**, 31, doi:10.1186/s13059-016-0893-4
650      (2016).
651  63  R Core Team. : A Language and Environment for Statistical Computing. *R Foundation for*
652      *Statistical Computing, Vienna, Austria, https://www.R-project.org*.

653

654

675    **Author contributions**

676    P.J.D.K, H.B. and T.F.M. designed experiments, P.J.D.K. and B.A.M.B. performed sBLISS experiments,

677    A.I. performed E.coli infection and immunostaining. Bioinformatics analysis were performed by P.J.D.K.

678    and H.B. Theoretical model of colibactin was built by F.B. and M.O. R.K. and L.A.A. provided and

679    analyzed WGS colorectal cancer data. The manuscript was written by P.J.D.K., H.B., F.B. and T.F.M.

25

680     **Competing interests**

681     The authors declare no competing interests

682     **Data and materials availability**

683     Input FASTA files and analysis scripts are available upon request. All other data is available in the main

684     text or supplementary materials.

685

26

686  **Figure Legends**

687



688

689  **Fig. 1. Identification of host DSB upon pks+ E.coli infection**

690  (A) Colibactin-producing E.coli infection causes γH2AX expression in Caco-2 cells.

691  (B) Experimental design for identification of positions of colibactin-induced DSBs with simplified BLISS
692  protocol.

693  (C) BLISS signal of etoposide-induced DSB shows increased counts compared to control condition.

694  (D) Heatmap indicating the log2 odds ratio of break enrichment in genomic region sets (FDR < 5%)
695  compared to the rest of the genome for *pks+* and *pks- E.coli* infected cells, etoposide treated and for
696  non-treated Caco-2 cells.

697

27

698

**Fig. 2. Colibactin damages DNA preferentially in specific AT-rich motifs**

(A) Enrichment of pentanucleotide sequences in close proximity to DSB positions (±3 nt) upon different treatments. Plots present pentanucleotide enrichment (log2 ratio of proportions of DSB at each motif between both conditions) of host breaks caused by pks+ E.coli infection in comparison to pks- E.coli infection and caused by pks- E.coli infection in comparison to breaks occurring in the non-treated (NT) cells, respectively.

(B) Consistency of an outstanding enrichment of AATTT and ATTTT and their reverse and complement sequences in colibactin induced DSBs in 4 independent biological replicates. Enrichment log2 ratios were standardized so that the highest log2 ratio of each experiment was taken to be 1 and the remaining values scaled accordingly. Enrichments are shown for 11 pentanucleotides with the highest standardized mean values. Each color refers to a different biological replicate.

(C) Preferred content of 5' and 3' dinucleotides next to colibactin's pentanucleotids motifs. For each of the motifs (AATTT and AAAAT) we first determined the log2 ratios for all 9nt sequences with the motif in the central 5nt. The 95% confidence interval was computed for each log2 ratio and the distribution of the lower bound of the interval plotted for each possible 2nt sequence at the 5' or 3' end of the central pentanucleotide.

(D) Top motif enriched in DSBs from pks+ E.coli infected cells compared to DSBs from pks- infected E.coli identified by Discriminative Regular Expression Motif Elicitation (DREME).

28
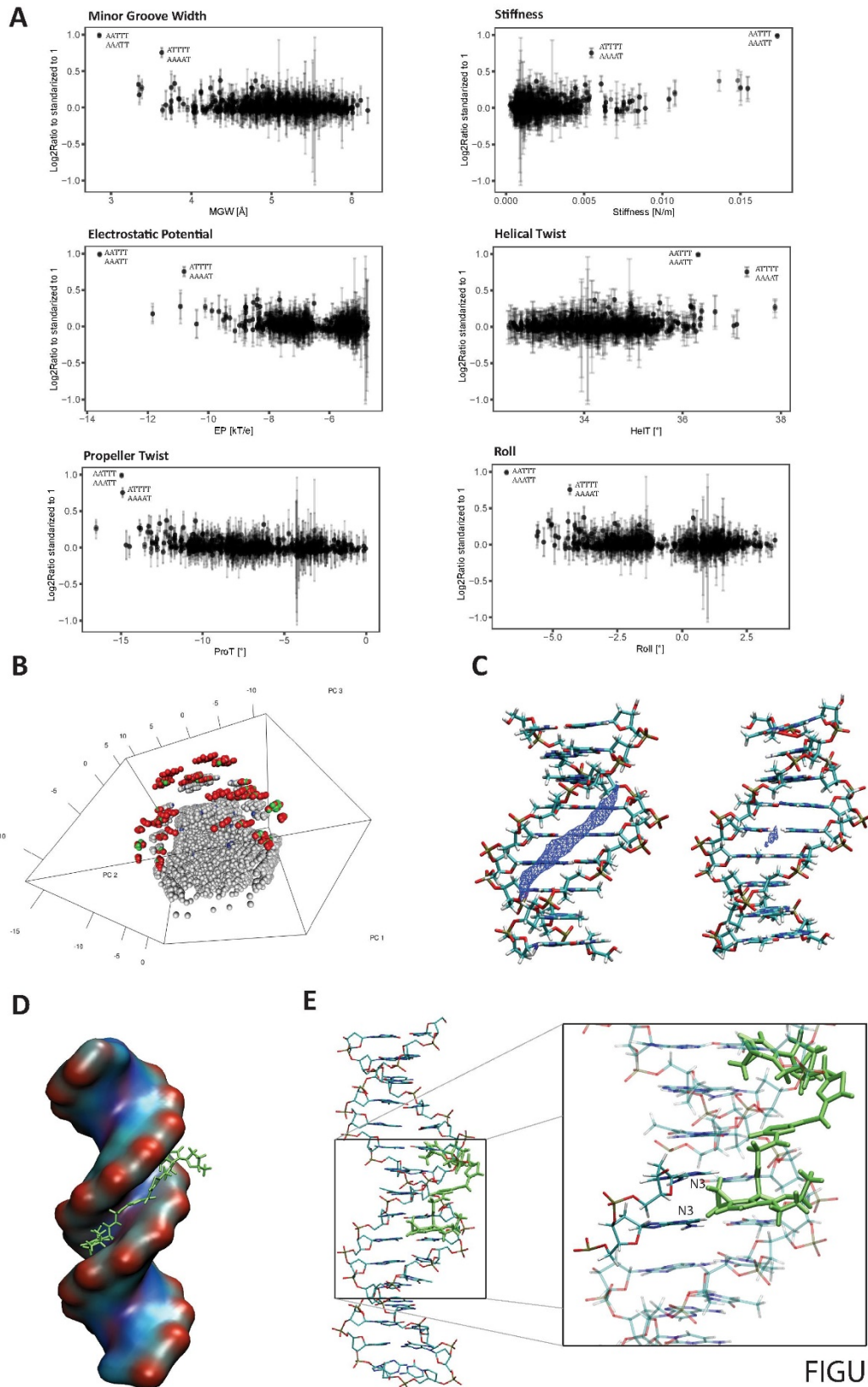
**Fig. 3. DSB caused by colibactin are associated with specific DNA shape pattern**

(A) Averaged minor groove width (MGW) predictions across all +/-100 bp contexts of DSBs identified by BLISS upon different treatments. The difference in averaged profiles of MGW for DSBs between the pks+ E.coli infection condition and all other treatments is enlarged and highlighted by a black arrow. As a control, MGW predictions of flanking sequences of 10,000 randomly chosen genomic positions are presented in the bottom right corner.

(B) MGW profiles of selected clusters obtained from k-means clustering of pks+ and pks- conditions. Landscape of cluster 1 for both conditions reflects the general pattern of MGW in close proximity to DSBs (note different y-axis scale). pks+ cluster 9 corresponds to AT-rich sequences across identified DSBs. Profiles of all parameters for every cluster can be found in Supplementary Fig 3A and B.

(C) Heatmap comparing averaged profiles of all identified clusters based on all predicted DNA shape parameters across pks+ and pks- infection conditions. Colors indicate individually Z-scored DNA shape characteristics. Each square in the heatmap refers to specific position from the break. Black arrows are marking exact DNA DSB position.Note that pks+ cluster 9 is unique for this treatment and shows extreme values centered at the DSB position.

29

FIGURE 4

733

**Fig. 4. Colibactin's binding motif corresponds to extreme DNA shape parameters values and extreme value of electrostatic potential**

(A) Correlation between pentanucleotide sequence enrichments (standardized to 1; for 4 biological replicates) for colibactin's activity and values of predicted DNA shape parameters. For MGW and EP values are calculated for each pentamer; for ProT intra-base pair parameter for the central base pair of each pentamer are calculated; for Roll and HelT the average of the two inter-base pair parameters, considering the two central base pair steps in each pentamer is calculated; for Stiffness average values, considering the two central base pair steps in each pentamer are calculated.

(B) 3D visualization of the first 3 principal components from predicted DNA shape values for the central 5nt of all possible 9nt motifs. Those 9nt motifs containing AAWWTT and/or showing strong enrichment are highlighted. Labels: red – AAWWTT motif with lower 95% confidence interval (CI) limit of log2ratio > 1.5; green – AAWWTT with lower CI limit < 1.5; blue – non-AAWWTT sequences with lower CI limit > 1.5, grey – other sequences (proportionally downsampled to approximately 35,000 sequences).

(C) Molecular Interaction Potential (MIP) using $Na^+$ as probe for 2 cases, on the left the most preferred 9nt DNA sequence for colibactin binding (CAAATTTTG) and on the right the least favorite (AACTTTGCA). The isosurfaces (in blue) for the two DNA sequences show different electrostatic potential (isovalue =-7 kcal $mol^{-1}$), correlating with the different minor groove conformations.

(D-E) Images of the theoretical docking of predicted colibactin structure into its preferred sequence motif (central sequence AAATTT), showing the insertion of the colibactin into the minor groove, with the double stranded DNA as surface (D) and showing the atomic details (E). Enlargement shows a zoomed-in image of the closeness of the cyclopropane to one of the N3 atom of the adenine, highlighting the possibility to alkylate the consequential base pair depending on the carbon involved in the alkylation.

31

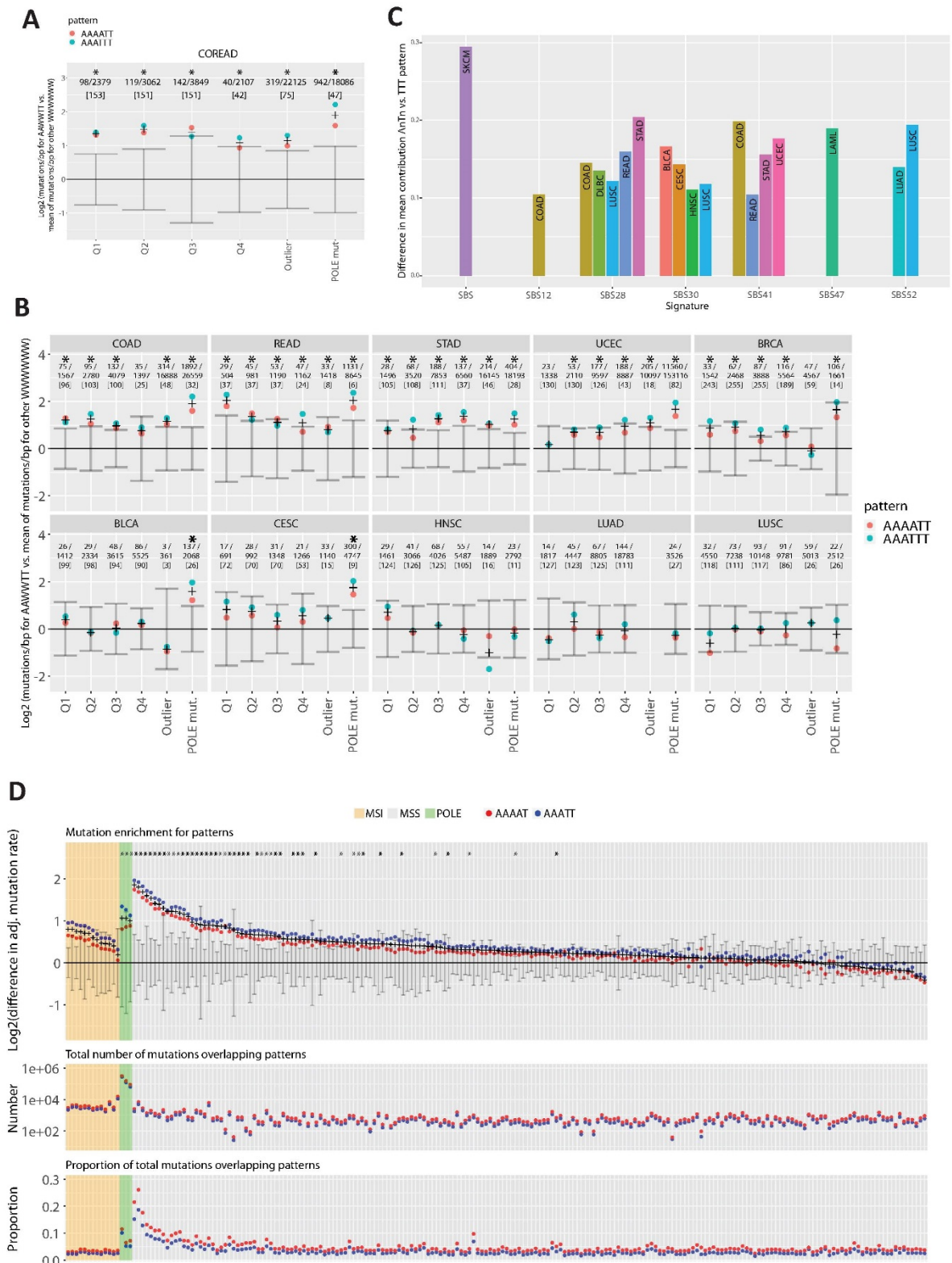FIGURE 5

**Fig. 5. Several cancers show enrichment of mutations at colbactin associated motifs**

A) Enrichment of single base change (SBS) mutations at colibactin-associated hexanucleotide motifs AAATTT/AAAATT in exome sequences from colorectal cancer cases [18]
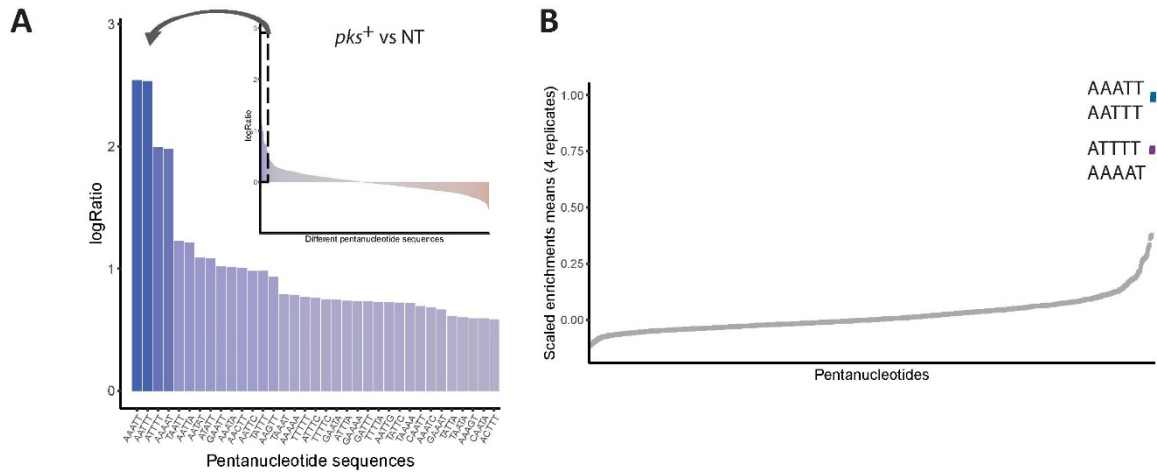
762    B) Enrichment of SBS mutations at colibactin associated hexanucleotide motifs AAATTT/AAAATT
763        in exome sequences from TCGA. Top row: cancer entities showing enrichment across all
764        subcohorts. Bottom row: cancer entities showing enrichment only for POLE mutated cases or
765        no enrichment at all. COAD-colon adenocarcinoma, READ-rectal adenocarcinoma, STAD-
766        stomach adenocarcinoma, UCEC-uterine corpus endometroid cancer, BRCA-breast cancer,
767        BLCA-bladder cancer, CESC-cervix squamous cell carcinoma, HNSC-head and neck squamous
768        cell cancer, LUAD-lung adenomcarcinoma, LUSC-lung squamous cell carcinoma
769    C) Signature detection rates for SBS mutations with contexts overlapping AATTT/ATTTT. Only
770        signatures with significant and positive differences in signature detection rates for contexts
771        overlapping AATTT/ATTTT compared to TTT are shown.
772    D)  Analysis of SBS mutations at colibactin associated pentanucleotide motifs AAATT/AAAAT  in
773        whole genome somatic mutation data from [19]. Top: Difference in log2(mutations/bp covered
774        by motif) between colibactin associated and all other WWWWW motifs. Middle: Total
775        mutation count at colibactin associated motifs: Bottom: Proportion of total mutations
776        ovlerapping colibactin associated motifs. MSS,MSI, and POLE, mutated cases.

777
778    (A), (B), (D) Stars denote significant difference (Mann-Whitney-U test $p < 0.05$ and FDR < 20%)
779        between colibactin associated motifs and all other motifs with the same A:T content and
780        length (A, B: hexanucleotide (HN) motif: AAATTT/AAAATT vs WWWWWW motifs,  C:
781        pentanucleotide (PN) motifs: AAATT/AAAAT vs. WWWWW motif). (A,B): First line is [number
782        of mutations overlapping AAWWTT motif] / [all mutations in cohort]. Third lines is number of
783        samples in cohort. Error bars describe the ± 2MAD intervals for mutation rate (mutations/bp
784        covered by motif) of WWWWW(W) motifs excluding colbactin associated pattern after
785        subtracting their mean. Dots represent the mutation rates for the two colibactin associated
786        PN or HN motifs after subtracting the mean of the WWWWW(W) motifs. Crosses are the
787        mean of the colibactin associated motifs.

1 **Supplementary Materials**

2

3 **Supplementary Figure Legends:**

4



5

6 **Fig. S1. Outstanding enrichment of AAATT and ATTTT and their reverse and complement sequences**
7 **in colibactin-induced DSBs.**

8 (A) Pentanucleotide sequences enriched (log2 ratio of proportions of DSB at each motif between
9 both conditions) at the DSB positions caused by colibactin-positive E.coli (pks+) in comparison to non-
10 treated (NT) cells.

11 (B) Scaled enrichment means of all 4 independent biological replicates for all possible
12 pentanucleotides (1,024) obtained by comparing nucleotide content of pks+ E.coli infection-induced
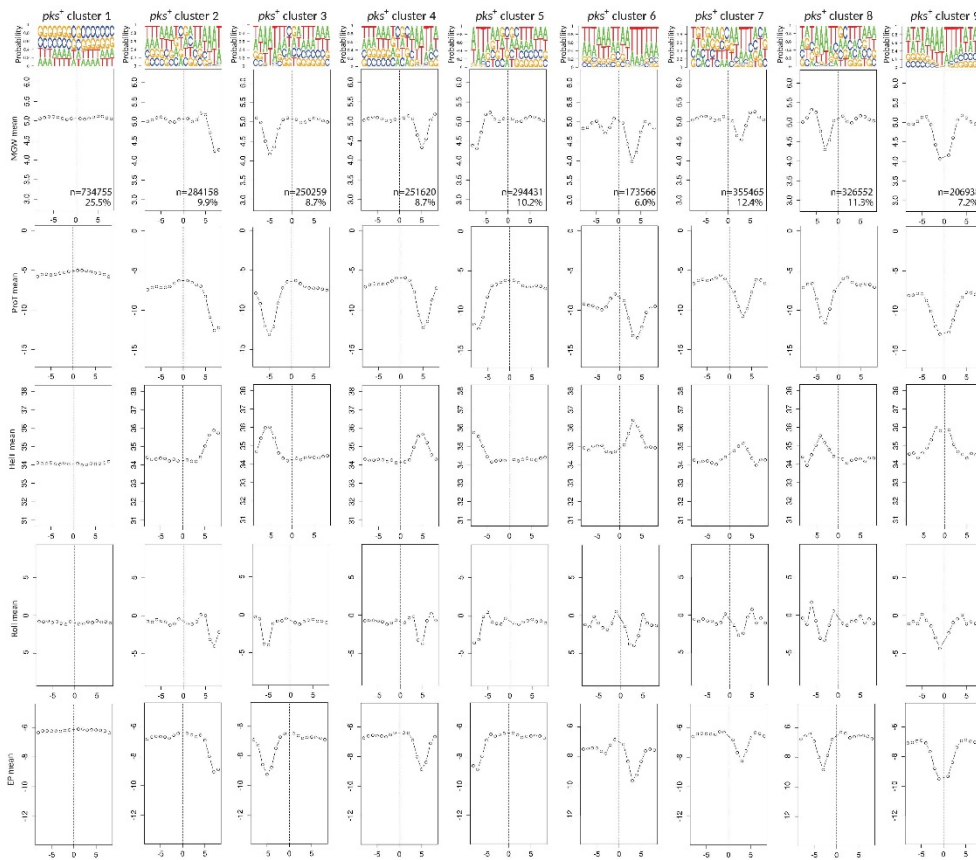13 breaks to the pks- E.coli-induced breaks.

1

FIGURE S2

14

**Fig. S2. Averaged values of DNA shape properties (HelT, ProT, Roll and EP) for sequences in proximity to identified DSBs of non-treated, etoposide-treated, pks- E.coli infected and pks+ E.coli infected cells.**
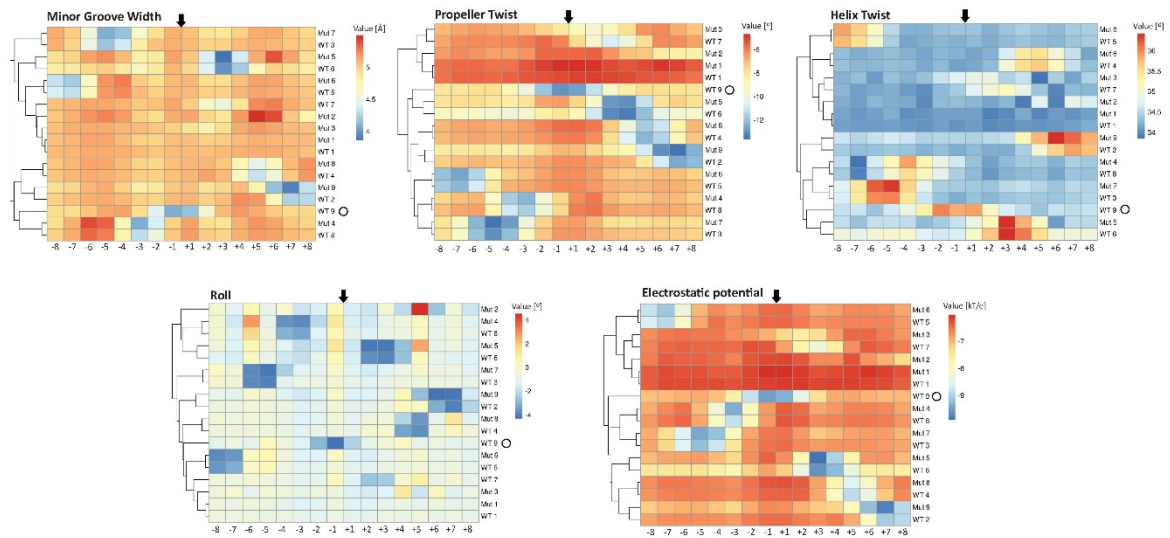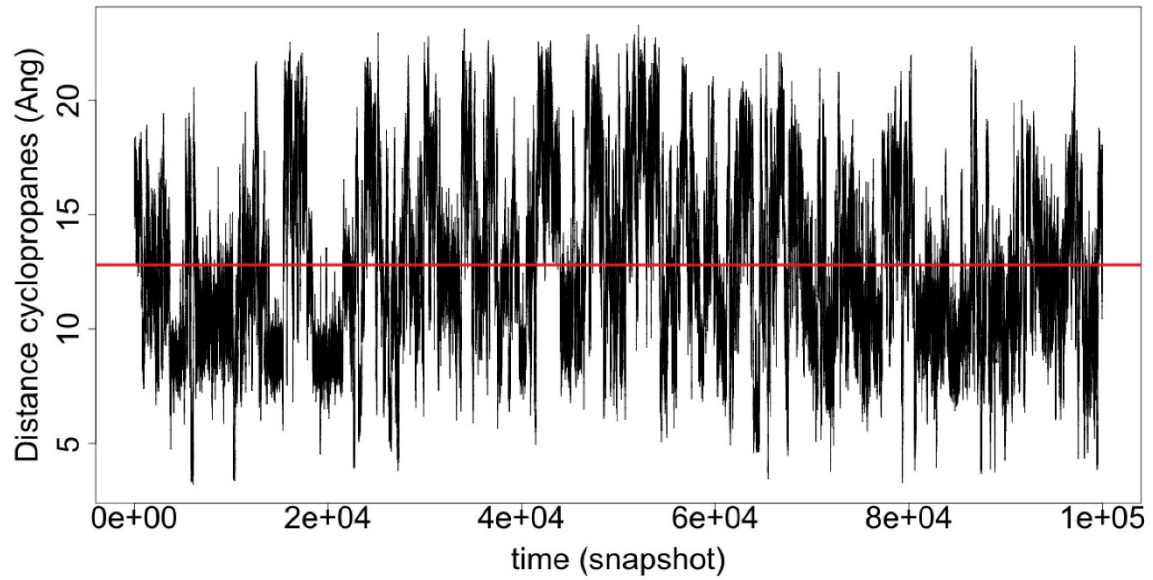
2

18



FIGURE S3a

19



FIGURE S3b

3

FIGURE S3c

**Fig. S3. K-means clustering of all predicted values of the DNA shape parameters.**

(A-B) DNA shape profiles of all clusters obtained from k-means clustering of pks+ and pks- conditions. (A) Clusters identified form the pks- dataset. (B) Clusters identified from the pks+ dataset. Above each cluster nucleotide probability for every position is presented.

(C) Heatmaps comparing averaged profiles of all identified clusters in pks+ and pks- conditions based on predicted DNA shape parameters , presented as individual comparisons of averaged profiles for each DNA shape parameter. Colors indicate absolute values for each DNA shape characteristics. Each square in the heatmap refers to specific position from the break. Black arrows are marking exact DNA DSB position.

4

31

**Fig. S4. Distance of the cyclopropanes (Å) along the MD simulation of the free colibactin in water. Red line identifies the average values of this distance (12.8Å).**

34

**Table S1** Summary of DSB enrichment at all 1024 pentanucleotide patterns across 4 replicates with associated predicted central DNA shape characteristics corresponding to Fig. 2A/B and Fig. 4A

37

5