

## 1 **Geographic and age-related variations in mutational processes in colorectal cancer**

2  
3 Marcos Díaz-Gay<sup>1,2,3,4,^</sup>, Wellington dos Santos<sup>5,^</sup>, Sarah Moody<sup>6,^</sup>, Mariya Kazachkova<sup>1,3,7</sup>, Ammal  
4 Abbasi<sup>1,2,3</sup>, Christopher D Steele<sup>1,2,3</sup>, Raviteja Vangara<sup>1,2,3</sup>, Sergey Senkin<sup>5</sup>, Jingwei Wang<sup>6</sup>, Stephen  
5 Fitzgerald<sup>6</sup>, Erik N Bergstrom<sup>1,2,3</sup>, Azhar Khandekar<sup>1,2,3,8</sup>, Burçak Otlu<sup>1,2,3,9</sup>, Behnoush Abedi-Ardekani<sup>5</sup>,  
6 Ana Carolina de Carvalho<sup>5</sup>, Thomas Cattiaux<sup>5</sup>, Ricardo Cortez Cardoso Penha<sup>5</sup>, Valérie Gaborieau<sup>5</sup>,  
7 Priscilia Chopard<sup>5</sup>, Christine Carreira<sup>10</sup>, Saamin Cheema<sup>6</sup>, Calli Latimer<sup>6</sup>, Jon W Teague<sup>6</sup>, Anush  
8 Mukeriya<sup>11</sup>, David Zaridze<sup>11</sup>, Riley Cox<sup>12</sup>, Monique Albert<sup>12,13</sup>, Larry Phouthavongsy<sup>12</sup>, Steven Gallinger<sup>14</sup>,  
9 Reza Malekzadeh<sup>15</sup>, Ahmadreza Niavarani<sup>15</sup>, Marko Miladinov<sup>16</sup>, Katarina Eric<sup>17</sup>, Sasa Milosavljevic<sup>18</sup>,  
10 Suleeporn Sangrajrang<sup>19</sup>, Maria Paula Curado<sup>20</sup>, Samuel Aguiar<sup>21</sup>, Rui Manuel Reis<sup>22,23</sup>, Monise Tadin  
11 Reis<sup>24</sup>, Luis Gustavo Romagnolo<sup>25</sup>, Denise Peixoto Guimarães<sup>26</sup>, Ivana Holcatova<sup>27,28</sup>, Jaroslav  
12 Kalvach<sup>29,30,31,32</sup>, Carlos Alberto Vaccaro<sup>33</sup>, Tamara Alejandra Piñero<sup>33</sup>, Beata Świątkowska<sup>34</sup>, Jolanta  
13 Lissowska<sup>35</sup>, Katarzyna Roszkowska-Purska<sup>36</sup>, Antonio Huertas-Salgado<sup>37</sup>, Tatsuhiro Shibata<sup>38,39</sup>, Satoshi  
14 Shiba<sup>39</sup>, Surasak Sangkhathat<sup>40,41,42</sup>, Taned Chitapanarux<sup>43</sup>, Gholamreza Roshandel<sup>44</sup>, Patricia Ashton-  
15 Prolla<sup>45,46</sup>, Daniel C Damin<sup>47</sup>, Francine Hehn de Oliveira<sup>48</sup>, Laura Humphreys<sup>6</sup>, Trevor D. Lawley<sup>49</sup>, Sandra  
16 Perdomo<sup>5</sup>, Michael R Stratton<sup>6</sup>, Paul Brennan<sup>5</sup>, and Ludmil B Alexandrov<sup>1,2,3,50,\*</sup>

17  
18 <sup>1</sup>Department of Cellular and Molecular Medicine, University of California San Diego, La Jolla, CA, USA

19 <sup>2</sup>Department of Bioengineering, University of California San Diego, La Jolla, CA, USA

20 <sup>3</sup>Moore's Cancer Center, University of California San Diego, La Jolla, CA, USA

21 <sup>4</sup>Digital Genomics Group, Structural Biology Program, Spanish National Cancer Research Center (CNIO),  
22 Madrid, Spain

23 <sup>5</sup>Genomic Epidemiology Branch, International Agency for Research on Cancer (IARC/WHO), Lyon,  
24 France

25 <sup>6</sup>Cancer, Ageing and Somatic Mutation, Wellcome Sanger Institute, Cambridge, UK

26 <sup>7</sup>Biomedical Sciences Graduate Program, University of California San Diego, La Jolla, CA, USA

27 <sup>8</sup>Division of Cancer Epidemiology and Genetics, National Cancer Institute, Bethesda, MD, USA

28 <sup>9</sup>Department of Health Informatics, Graduate School of Informatics, Middle East Technical University,  
29 Ankara, Turkey

30 <sup>10</sup>Evidence Synthesis and Classification Branch, International Agency for Research on Cancer  
31 (IARC/WHO), Lyon, France

32 <sup>11</sup>Clinical Epidemiology, N.N. Blokhin National Medical Research Centre of Oncology, Moscow, Russia

33 <sup>12</sup>Ontario Tumour Bank, Ontario Institute for Cancer Research, Toronto, ON, Canada

34 <sup>13</sup>Centre for Biodiversity Genomics, University of Guelph, Guelph, ON, Canada

35 <sup>14</sup>Lunenfeld-Tanenbaum Research Institute, Sinai Health System, Toronto, ON, Canada

36 <sup>15</sup>Digestive Oncology Research Center, Digestive Disease Research Institute, Tehran University of Medical  
37 Sciences, Tehran, Iran

38 <sup>16</sup>Clinic for Digestive Surgery - First Surgical Clinic, University Clinical Centre of Serbia, Belgrade, Serbia

39 <sup>17</sup>Department of Pathology, University Clinical Centre of Serbia, Belgrade, Serbia

40 <sup>18</sup>International Organization for Cancer Prevention and Research, Belgrade, Serbia

41 <sup>19</sup>National Cancer Institute, Bangkok, Thailand

42 <sup>20</sup>Department of Epidemiology, A.C. Camargo Cancer Center, Sao Paulo, Brazil

43 <sup>21</sup>Colon Cancer Reference Center, A.C. Camargo Cancer Center, Sao Paulo, Brazil

44 <sup>22</sup>Molecular Oncology Research Center, Barretos Cancer Hospital, Barretos, Brazil

45 <sup>23</sup>Life and Health Sciences Research Institute (ICVS), School of Medicine, Minho University, Braga,  
46 Portugal

47 <sup>24</sup>Department of Pathology, Barretos Cancer Hospital, Barretos, Brazil

48 <sup>25</sup>Department of Colorectal Oncology Surgery, Barretos Cancer Hospital, Barretos, Brazil

49 <sup>26</sup>Department of Endoscopy, Barretos Cancer Hospital, Barretos, Brazil

50 <sup>27</sup>Institute of Public Health & Preventive Medicine, 2<sup>nd</sup> Faculty of Medicine, Charles University, Prague,  
51 Czech Republic

52 <sup>28</sup>Department of Oncology, 2<sup>nd</sup> Faculty of Medicine, Charles University and Motol University Hospital,  
53 Prague, Czech Republic

54 <sup>29</sup>Surgery Department, 2<sup>nd</sup> Faculty of Medicine, Charles University and Central Military Hospital, Prague,  
55 Czech Republic

56 <sup>30</sup>2<sup>nd</sup> Faculty of Medicine, Charles University and Motol University Hospital, Prague, Czech Republic

57 <sup>31</sup>Institute of Animal Physiology and Genetics Czech Academy of Science, Libechov, Czech Republic

58 <sup>32</sup>Clinical Center ISCARE, Prague, Czech Republic

59 <sup>33</sup>Instituto de Medicina Traslacional e Ingeniería Biomédica (IMTIB)- CONICET- Universidad Hospital  
60 Italiano de Buenos Aires (UHIBA) y Hospital Italiano de Buenos Aires (HIBA), Buenos Aires, Argentina

61 <sup>34</sup>Department of Environmental Epidemiology, Nofer Institute of Occupational Medicine, Łódź, Poland

62 <sup>35</sup>The Maria Sklodowska-Curie National Research Institute of Oncology, Warsaw, Poland

63 <sup>36</sup>Department of Pathology, The Maria Sklodowska-Curie National Research Institute of Oncology,  
64 Warsaw, Poland

65 <sup>37</sup>Oncological pathology group, Terry Fox National Tumor Bank (Banco Nacional de Tumores Terry Fox),  
66 National Cancer Institute, Bogotá, Colombia

67 <sup>38</sup>Laboratory of Molecular Medicine, The Institute of Medical Science, The University of Tokyo, Minato-  
68 ku, Japan

69 <sup>39</sup>Division of Cancer Genomics, National Cancer Center Research Institute, Chuo-ku, Japan

70 <sup>40</sup>Translational Medicine Research Center, Faculty of Medicine, Prince of Songkla University, Hat Yai,  
71 Thailand

72 <sup>41</sup>Department of Biomedical Sciences and Biomedical Engineering, Faculty of Medicine, Prince of Songkla  
73 University, Hat Yai, Thailand

74 <sup>42</sup>Department of Surgery, Faculty of Medicine, Prince of Songkla University, Hat Yai, Thailand

75 <sup>43</sup>Department of Internal Medicine, Faculty of Medicine, Chiang Mai University, Chiang Mai, Thailand

76 <sup>44</sup>Golestan Research Center of Gastroenterology and Hepatology, Golestan University of Medical Sciences,  
77 Gorgan, Iran

78 <sup>45</sup>Department of Genetics, Universidade Federal do Rio Grande do Sul (UFRGS), Porto Alegre, Brazil

79 <sup>46</sup>Medical Genetics Service, Hospital de Clínicas de Porto Alegre (HCPA), Porto Alegre, Rio Grande do  
80 Sul, Brazil

81 <sup>47</sup>Department of Surgery, Division of Colorectal Surgery, Hospital de Clínicas de Porto Alegre (HCPA),  
82 Porto Alegre, Rio Grande do Sul, Brazil

83 <sup>48</sup>Department of Pathology, Anatomic Pathology, Hospital de Clínicas de Porto Alegre (HCPA), Porto  
84 Alegre, Rio Grande do Sul, Brazil

85 <sup>49</sup>Parasites and Microbes, Wellcome Sanger Institute, Cambridge, UK

86 <sup>50</sup>Sanford Stem Cell Institute, University of California San Diego, La Jolla, CA, USA

87

88 ^These authors contributed equally

89 \*Correspondence should be addressed to [L2alexandrov@health.ucsd.edu](mailto:L2alexandrov@health.ucsd.edu)

90

91 **ABSTRACT**

92 Colorectal cancer incidence rates vary geographically and have changed over time. Notably, in the  
93 past two decades, the incidence of early-onset colorectal cancer, affecting individuals under the  
94 age of 50 years, has doubled in many countries. The reasons for this increase are unknown. Here,  
95 we investigate whether mutational processes contribute to geographic and age-related differences  
96 by examining 981 colorectal cancer genomes from 11 countries. No major differences were found  
97 in microsatellite unstable cancers, but variations in mutation burden and signatures were observed  
98 in the 802 microsatellite-stable cases. Multiple signatures, most with unknown etiologies,  
99 exhibited varying prevalence in Argentina, Brazil, Colombia, Russia, and Thailand, indicating  
100 geographically diverse levels of mutagenic exposure. Signatures SBS88 and ID18, caused by the  
101 bacteria-produced mutagen colibactin, had higher mutation loads in countries with higher  
102 colorectal cancer incidence rates. SBS88 and ID18 were also enriched in early-onset colorectal  
103 cancers, being 3.3 times more common in individuals diagnosed before age 40 than in those over  
104 70, and were imprinted early during colorectal cancer development. Colibactin exposure was  
105 further linked to *APC* driver mutations, with ID18 responsible for about 25% of *APC* driver indels  
106 in colibactin-positive cases. This study reveals geographic and age-related variations in colorectal  
107 cancer mutational processes, and suggests that early-life mutagenic exposure to colibactin-  
108 producing bacteria may contribute to the rising incidence of early-onset colorectal cancer.

## 109 INTRODUCTION

110 The age-standardized incidence rates (ASR) for most adult cancers vary across different  
111 geographic locations and can change over time<sup>1</sup>. Despite extensive epidemiological research, the  
112 underlying causes for many of these variations remain unclear. However, they are suspected to be  
113 due to exogenous environmental or lifestyle carcinogenic exposures, which are, in principle,  
114 preventable<sup>2</sup>. Many well-known exogenous carcinogens are also mutagens<sup>3,4</sup>, which can imprint  
115 characteristic patterns of somatic mutations in the genome, known as mutational signatures.  
116 Therefore, a complementary approach to conventional epidemiology for investigating unknown  
117 causes of cancer is the characterization of mutational signatures in the genomes of cancer and  
118 normal cells<sup>5-7</sup>. The *Mutographs* Cancer Grand Challenge project<sup>8</sup> has implemented this strategy  
119 of “mutational epidemiology” by sequencing cancers from geographic areas of differing incidence  
120 rates, using mutational signature analysis to elucidate the mutational processes that have been  
121 operative, with results thus far from cancers of the esophagus<sup>6</sup>, kidney<sup>5</sup>, and head and neck<sup>9</sup>.

122  
123 Colorectal cancer incidence rates differ markedly by geographic location and have changed  
124 substantially in some countries over the last 70 years<sup>10</sup>. For instance, the ASR for colorectal cancer  
125 in North America and in most European countries peaked in the 1980s and 1990s and have been  
126 declining since, whereas countries in East Asia such as Japan and South Korea have been steadily  
127 increasing over the past seven decades<sup>1</sup>. Moreover, in the past 20 years there has been a notable  
128 global increase in the incidence of early-onset colorectal cancer<sup>10,11</sup>, typically defined as colorectal  
129 cancer in adults under 50 years of age. This was first reported in the United States<sup>12</sup> and  
130 subsequently observed in Australia, Canada, Japan<sup>13</sup> and multiple European countries<sup>14</sup>. Although  
131 epidemiological studies have identified multiple risk factors for colorectal cancer, specific risk

132 factors for early-onset colorectal cancer remain largely unidentified, with the exception of family  
133 history and hereditary predisposition. The latter is predominantly attributable to Lynch syndrome,  
134 which is characterized by DNA mismatch repair deficient cancers of the proximal colon<sup>15,16</sup> and,  
135 therefore, is unlikely to be implicated in the recent increase in early-onset colorectal cancer, which  
136 is mainly enriched in sporadic, DNA mismatch repair proficient cancers affecting the distal colon  
137 and rectum<sup>17,18</sup>.

138  
139 Previous colorectal cancer whole-genome sequencing studies have largely focused on cases from  
140 North America and Europe including USA<sup>19,20</sup>, UK<sup>21,22</sup>, Netherlands<sup>23-26</sup>, and Sweden<sup>27</sup> and  
141 incorporated limited numbers of early-onset cases<sup>19,21,22,26,27</sup>. Here, we examine colorectal cancer  
142 genomes from 11 countries on four continents to investigate whether variation in mutational  
143 processes contributes to geographic and age-related differences in incidence rates.

144

## 145 RESULTS

### 146 Study design

147 981 colorectal cancers (45.7% female) were collected from intermediate-incidence countries with  
148 ASRs of 13-20/100,000 people (Iran, Thailand, Colombia, Brazil) and high-incidence countries  
149 with ASRs >24 (Argentina, Canada, Russia, Serbia, Czech Republic, Poland, Japan), including the  
150 highest ASR of 37 in Japan<sup>1</sup> (**Fig. 1a; Supplementary Table 1**). Of the 981 cases, 320 were from  
151 the proximal colon, 333 from the distal colon, 326 from the rectum, and 2 from unspecified subsites  
152 (**Fig. 1b**). There were 132 early-onset cases, which were 1.88-fold enriched in the distal colon and  
153 rectum compared to the proximal colon ( $p=0.006$ ). All cancers and their matched normal samples  
154 underwent whole genome sequencing, achieving a median coverage of 53-fold and 27-fold,  
155 respectively.

156

### 157 Mutation burden and molecular classification

158 The 981 colorectal cancers were divided into known molecular subtypes based on their somatic  
159 mutation burdens and profiles. Consistent with prior studies<sup>20,28</sup>, two main subtypes were  
160 identified: DNA mismatch repair proficient cancers, also known as microsatellite stable (MSS),  
161 and DNA mismatch repair deficient cancers, often referred to as tumors showing microsatellite  
162 instability (MSI). MSS samples ( $n=802$ , 81.8%; **Fig. 1c**) were characterized by a lower burden of  
163 single base substitutions (SBS; median: 12,054) and small insertions and deletions (ID; median:  
164 1,451), and a higher burden of large-scale genomic aberrations (median: 53.5% of genome altered).  
165 In contrast, MSI samples ( $n=153$ , 15.6%) exhibited higher SBS and ID burdens (median: 95,426  
166 and 125,100, respectively) with limited genomic aberrations (median: 7.0%). As expected, the

167 average mutational profiles of MSS and MSI colorectal tumors were different (**Extended Data**  
168 **Fig. 1a-b**).

169  
170 MSI samples were predominantly found in the proximal colon (OR=12.2,  $p=3.8 \times 10^{-27}$ ) and were  
171 more common in early-onset cases (OR=2.6,  $p=0.001$ ). Notably, 31/153 MSI cases (20.3%),  
172 including 13/28 MSI early-onset cases (46.4%), carried germline pathogenic variants in DNA  
173 mismatch repair genes consistent with Lynch syndrome (**Supplementary Table 2**). After  
174 excluding all cases attributed to Lynch syndrome, there was no enrichment of MSI cancers in  
175 early-onset cases ( $p>0.05$ ). Deficiencies of other DNA repair mechanisms were observed in 24/981  
176 cancers (2.4%), including ultra-hypermutated cases with mutations in *POLE* ( $n=10$ , 1.0%) and  
177 *POLD1* polymerases ( $n=3$ , 0.3%), homologous recombination deficient (HRD) cases ( $n=7$ , 0.7%),  
178 and cases with mutations in the base excision repair genes *MUTYH* ( $n=1$ , 0.1%), *NTHL1* ( $n=2$ ,  
179 0.2%), and *OGGI* ( $n=1$ , 0.1%) (**Methods; Supplementary Table 3-4; Supplementary Fig. 1-3**).

180  
181 The mutational catalogues of DNA repair deficient cancers are dominated by somatic mutations  
182 resulting from the failed repair process, rendering it difficult to characterize mutational processes  
183 unrelated to this failure<sup>29</sup>. To enable investigation of the latter, we therefore focused the main  
184 analyses on DNA repair proficient colorectal cancers, while reporting DNA repair deficient cases  
185 in the **Supplementary Note**. Two cases treated with chemotherapy for prior cancers were also  
186 excluded as their mutation profiles were dominated by the mutational signatures of chemotherapy  
187 agents<sup>19,30</sup> (**Supplementary Fig. 4**). The remaining cohort consisted of 802 treatment-naïve DNA  
188 repair proficient colorectal cancers, including 97 early-onset cases.

189

190 After adjustment for sex, country, tumor subsite, and tumor purity (**Methods**), early-onset cancers  
191 showed reduced burdens of SBS (fold-change [FC]=0.92,  $p=0.045$ ) and ID (FC=0.90,  $p=0.018$ ;  
192 **Fig. 1d**) but not of doublet base substitutions (DBS), copy number alterations (CN), or structural  
193 variants (SV) when compared to late-onset cases ( $p>0.05$ ). Nevertheless, the average mutation  
194 spectra of early-onset and late-onset cancers were remarkably similar for all types of somatic  
195 mutations (cosine similarity $>0.97$ ; **Fig. 1e-g**; **Extended Data Fig. 1c-d**). Mutation burden also  
196 varied substantially for specific countries when compared to all others, including Canada (lower  
197 SBS and ID burdens), Poland (higher SBS and DBS), Japan (lower SBS, ID, DBS), Iran (lower  
198 ID), and Brazil (higher ID and CN; **Extended Data Fig. 2**). However, mutation profiles were  
199 generally consistent across all countries (**Extended Data Fig. 3**).

200

## 201 **Repertoire of mutational signatures**

202 A total of 16 SBS, 10 ID, 4 DBS, 6 CN, and 6 SV *de novo* mutational signatures were extracted  
203 from the 802 MSS colorectal cancers and subsequently decomposed into a combination of  
204 previously reported reference signatures and potential novel signatures (**Supplementary Note**;  
205 **Supplementary Tables 5-15**). The 16 *de novo* SBS signatures encompassed 15 COSMICv3.4  
206 signatures (**Extended Data Fig. 4a**; **Supplementary Table 10**), including those previously  
207 associated with clock-like mutational processes (SBS1, SBS5)<sup>31</sup>, APOBEC deamination (SBS2,  
208 SBS13)<sup>31</sup>, deficient homologous recombination (SBS3)<sup>31</sup>, reactive oxygen species (SBS18)<sup>32</sup>,  
209 exposure to the mutagenic agent colibactin synthesized by *Escherichia coli* and other microbes  
210 carrying a ~40kb polyketide synthase (*pks*) pathogenicity island (SBS88)<sup>33,34</sup>, and mutational  
211 processes of unknown causes (SBS8, SBS17a/b, SBS34, SBS40a, SBS89, SBS93,  
212 SBS94)<sup>5,19,32,34,35</sup>. Three previously described signatures of unknown origin<sup>21</sup> (SBS\_F, SBS\_H,



213 SBS\_M; **Extended Data Fig. 4b**) and a novel signature (SBS\_O; **Extended Data Fig. 4c**) were  
214 also detected. SBS\_O corresponds to a refined version of a previously reported signature of  
215 unknown etiology (SBS41; **Methods**)<sup>19</sup>. With respect to ID, DBS, CN, and SV, most *de novo*  
216 extracted mutational signatures were highly similar to, or directly reconstructed by, COSMICv3.4  
217 reference signatures (**Extended Data Fig. 4d-f** and **5**; **Supplementary Table 10**) with the  
218 exception of an ID signature (ID\_J), characterized by deletions of isolated Ts and insertions of Ts  
219 in long repetitive regions resembling a previously reported signature<sup>34</sup> (**Extended Data Fig. 4e**),  
220 and three novel signatures from large mutational events (CN\_F, SV\_B, SV\_D; **Extended Data**  
221 **Fig. 5b&d**), which were extracted due to the extended contexts used in our signature analysis  
222 (**Methods**).

223

#### 224 **Geographic variation in mutational signatures**

225 Despite the similar mutation profiles across countries (**Extended Data Fig. 3**), several signatures  
226 exhibited varying prevalence when comparing one country to all others (**Fig. 2a**; **Supplementary**  
227 **Fig. 5**; **Supplementary Table 16**). Notably, SBS89 (OR=28.0,  $q=0.001$ ), DBS8 (OR=8.9,  
228  $q=3.2\times 10^{-4}$ ), and the novel ID\_J (OR=9.6,  $q=6.2\times 10^{-5}$ ) were at elevated frequencies in Argentina  
229 when compared to all other countries (**Fig. 2b**). Signatures SBS89, DBS8, and ID\_J also showed  
230 a strong tendency to co-occur ( $p<1.7\times 10^{-11}$ ) suggesting they may arise from the same underlying  
231 mutational process. In Colombia (**Fig. 2c**), higher frequencies were observed for SBS94 (OR=19.7,  
232  $q=3.2\times 10^{-5}$ ), the novel SBS\_F (OR=10.7,  $q=2.0\times 10^{-4}$ ), and DBS6 (OR=12.5,  $q=0.028$ ) when  
233 compared to all other countries, with evidence of co-occurrence of SBS94 with SBS\_F ( $p=0.017$ )  
234 and DBS6 ( $p=1.9\times 10^{-4}$ ). Enrichments were also found for SBS2 (OR=2.0,  $q=0.041$ ) and SBS\_H  
235 (OR=2.3,  $q=0.001$ ) in Russia and CN\_F (OR=3.5,  $q=3.9\times 10^{-4}$ ) in Brazil, whereas depletions were

236 identified for DBS2 in Thailand (OR=0.38,  $q=0.008$ ) and for DBS4 in Colombia (OR=0.06,  
237  $q=0.034$ ; **Fig. 2a**). Overall, the results indicate international differences in the prevalence of certain  
238 mutational processes involved in colorectal cancer development.

239  
240 To explore the broader epidemiological implications of international variation in mutational  
241 processes, as previously done for kidney cancer<sup>5</sup>, we evaluated the relationships between ASR and  
242 mutational signatures (**Fig. 2d; Supplementary Table 17**). Independent of covariates, colibactin-  
243 induced mutational signatures, SBS88 and ID18, as well as clock-like signature SBS1 and novel  
244 signature SBS\_H, associated with an increasing rate of ASR for colorectal cancer, whereas novel  
245 signature CN\_F associated with a reduced ASR rate ( $q<0.05$ ; **Fig. 2d-e; Extended Data Fig. 6a**).  
246 For SBS88 and ID18, the association was linked with the ASR for rectal cancer ( $q=0.088$  and  
247  $q=0.008$ ; **Fig. 2f; Supplementary Table 18**). In contrast, for SBS1, SBS\_H, and CN\_F the  
248 association was particularly strong for the ASR of colon cancer ( $q=0.009$ ,  $q=0.015$ , and  $q=0.057$ ;  
249 **Extended Data Fig. 6b**). Colibactin-associated signatures were also found elevated in patients  
250 from countries with high ASR rates for early-onset colorectal cancer (**Extended Data Fig. 6c**).

251  
252 **Colibactin induced mutational signatures are enriched in early-onset colorectal cancer**

253 In addition to examining the global distribution of mutational signatures, the substantial number  
254 of early-onset colorectal cancer cases enabled evaluating the association between mutational  
255 signatures and age at diagnosis. Although the average mutation profiles of early-onset and late-  
256 onset colorectal cancer cases were similar (**Fig. 1e-g**), the prevalence of some mutational  
257 signatures was associated with the age of diagnosis, independently of country of origin (**Fig. 3a**;  
258 **Supplementary Table 19**), genetic ancestry or ethnicity (**Supplementary Fig. 6-8**). As expected,

259 late-onset cases showed enrichment in signatures known to accumulate linearly with age in normal  
260 colorectal crypts<sup>36</sup>, including SBS1, SBS5, ID1, and ID2 (**Fig. 3a-b**). Unknown etiology indel  
261 signatures ID4, ID9, and ID10 also showed associations with late-onset cases (**Fig. 3a-b**).

262  
263 By contrast, enrichment in early-onset cancers was observed for colibactin-induced signatures.  
264 Signatures SBS88 and ID18 were 2.5 and 4 times more common, respectively, in colorectal  
265 cancers diagnosed below than above the age of 50 ( $q=0.006$  and  $q=3.7\times 10^{-7}$ , respectively; **Fig. 3a-**  
266 **b**). The primary associations of early-onset cases with SBS88 and ID18 were further supported by  
267 the successive decline in the prevalence of these signatures with increasing age of diagnosis ( $p$ -  
268  $trend=1.3\times 10^{-4}$  and  $p$ - $trend=2.0\times 10^{-7}$ , respectively; **Fig. 3c**; **Supplementary Table 20**). A similar  
269 effect was observed using a complementary motif enrichment analysis for detecting SBS88,  
270 similarly to a recent study<sup>26</sup> ( $p$ - $trend=1.0\times 10^{-7}$ ; **Extended Data Fig. 7a-b**). On the basis of the  
271 strong co-occurrence of SBS88 and ID18 ( $p=7.4\times 10^{-63}$ ), as well as previous functional<sup>33</sup> and  
272 population studies<sup>22,26,34</sup>, we defined exposure to colibactin by the presence of either SBS88 or  
273 ID18. Colibactin exposure was found in 21.1% of all colorectal cancers (169/802) and was  
274 associated with earlier age of onset (median age: 62 vs. 67,  $p=1.6\times 10^{-8}$ ; **Fig. 3d**), an effect more  
275 evident in the distal colon (median age: 57 vs. 66,  $q=5.2\times 10^{-7}$ ) and rectum (median age: 63 vs. 66,  
276  $q=0.025$ ; **Fig. 3e**). Overall, colibactin exposure had a strong inverse correlation with age, being  
277 3.3 times more common in colorectal cancers diagnosed in individuals younger than 40 compared  
278 to those over 70 ( $p$ - $trend=2.7\times 10^{-7}$ , **Extended Data Figure 7c**).

279  
280 Signatures of unknown etiology SBS\_M and ID14 (**Fig. 3a-c**) were also enriched in early-onset  
281 cases, and SBS89 similarly exhibited a higher prevalence in younger individuals (5.8 times more

282 prevalent in early-onset compared to late-onset patients with  $p$ -trend=0.047), albeit based on a  
283 very small number of cancers harboring the signature (9/802, 1.1%; **Fig. 3c**). Interestingly, SBS\_M  
284 showed an elevation in distal colon and rectum tumors compared to proximal colon similar to the  
285 one observed in colibactin-associated signatures SBS88 and ID18, previously reported<sup>22</sup>  
286 (**Supplementary Fig. 9**).

287

### 288 **Colibactin mutagenesis is an early event in colorectal carcinogenesis**

289 To time the imprinting of SBS88 and ID18, mutations were categorized as early clonal, late clonal,  
290 or subclonal during the development of each cancer and the contribution of each mutational  
291 signature to each category was determined (**Methods**). SBS88 and ID18 were both enriched in  
292 early clonal compared to late clonal mutations ( $q=4.2\times 10^{-4}$  and  $q=6.1\times 10^{-5}$ ; **Fig. 4a**), as well as a  
293 similar trend in clonal compared to subclonal mutations ( $q=0.138$  and  $q=0.058$ ; **Extended Data**  
294 **Fig. 8a**), consistent with the presence of these mutational signatures in normal colorectal  
295 epithelium<sup>34</sup>. This enrichment in earlier evolutionary stages was similar to the one observed for  
296 other well-known clock-like signatures like SBS1, SBS5, or ID1 (**Fig. 4a-b**), as previously shown  
297 in tumors<sup>37,38</sup> and normal tissues<sup>34</sup>, and in contrast to signatures known to preferentially generate  
298 late clonal and subclonal mutations, such as SBS17a/b<sup>38</sup>. Interestingly, the enrichment of colibactin  
299 signatures in early clonal mutations was observed for both early-onset ( $q=0.004$  for SBS88 and  
300  $q=2.0\times 10^{-4}$  for ID18) and late-onset colorectal cancer cases ( $q=0.020$  and  $q=0.024$ ; **Extended**  
301 **Data Fig. 8b**).

302

303 Since colibactin is produced by bacteria carrying the *pks* pathogenicity island, we investigated  
304 whether colorectal cancer cases with SBS88 or ID18 harbored *pks*+ bacteria based on sequencing

305 reads from the cancer sample that did not map to the human genome but mapped to the *pks* locus  
306 (**Methods**). Consistent with a prior observation<sup>39</sup>, there was no association between the presence  
307 of SBS88 or ID18 and that of *pks*<sup>+</sup> bacteria (**Fig. 4b; Extended Data Fig. 9**). Similarly, no  
308 microbiome association was observed for the other signatures enriched in early-onset colorectal  
309 cancers (**Supplementary Note**). Moreover, we observed a younger age of diagnosis for cases with  
310 SBS88 or ID18 but without an identified *pks*<sup>+</sup> bacteria ( $p=1.3\times 10^{-7}$ ; **Fig. 4c-d**). While the reasons  
311 are unclear, one likely explanation is the imprinting of SBS88 and ID18 on the colorectal  
312 epithelium during an early period of life when *pks*<sup>+</sup> bacteria were present, followed by the natural  
313 plasticity of the microbiome over subsequent decades, leading to the loss or gain of *pks*<sup>+</sup> bacteria.  
314

### 315 **Colibactin exposure and driver mutations**

316 Using the IntOGen framework<sup>40</sup>, 46 genes under positive selection were identified, with eight  
317 mutated in more than 10% of cancers: *APC*, *TP53*, *KRAS*, *FBXW7*, *SMAD4*, *PIK3CA*, *TCFL2*, and  
318 *SOX9* (**Fig. 5a; Supplementary Table 21**). Forty-three of the 46 genes have been previously  
319 reported as colorectal cancer driver genes<sup>22,40</sup>, two in other cancer types (*MED12*, *NCOR1*)<sup>40</sup>, and  
320 a putative novel colorectal cancer driver gene, *CCR4*, was identified with mutations indicating  
321 inactivation of the encoded protein. Mutations affecting these 46 cancer driver genes were  
322 annotated as driver mutations using a multi-step process based on the mutation type and the mode  
323 of action of the gene (**Methods**). An elevation in the total number of driver mutations was observed  
324 in late-onset compared to early-onset cases ( $FC=1.21$ ,  $p=5.4\times 10^{-5}$ ; **Fig. 5b**). In addition, an  
325 enrichment in *APC* driver mutation carriers was also found for late-onset cases ( $OR=2.7$ ,  $q=0.027$ ;  
326 **Fig. 5c-d; Supplementary Table 22**), as previously reported<sup>41</sup>, whereas no hotspot driver mutation  
327 (defined as those affecting the same genomic position in at least 10 cases) was associated with age

328 of onset ( $q>0.05$ ; **Supplementary Table 23**). No statistically significant differences across  
329 countries were found for driver mutations within cancer driver genes or for hotspot driver  
330 mutations ( $q>0.05$ ; **Supplementary Tables 24-25**).

331  
332 The contributions of SBS88 and ID18 to driver mutations were assessed using probabilistic  
333 assignment of signatures to individual mutations<sup>42</sup>. SBS88 accounted for 64.3% of the colibactin-  
334 induced<sup>43</sup> *APC* splicing variant c.835-8A>G in colibactin-exposed samples, compared to only  
335 3.9% and 3.8% of driver substitutions in *APC* or other cancer genes (**Fig. 5e**). Similarly, ID18  
336 accounted for 25.3% of *APC* driver indels and 16.9% of other driver indels in colibactin-exposed  
337 cases (**Fig. 5f**). Overall, SBS88 and ID18 accounted for 8.3% of all SBS and ID driver mutations,  
338 and 15.5% of all *APC* driver mutations in colibactin positive cancers. Nevertheless, no differences  
339 were observed between early-onset and late-onset colibactin positive colorectal cancer in the  
340 proportion of driver mutations assigned to specific mutational signatures (**Extended Data Fig.**  
341 **10a-b**). In addition, a prior study observed that SBS88 is also responsible for mutations in  
342 chromatin modifier genes<sup>39</sup>, and we were able to validate this as well as show a similar effect for  
343 the colibactin-associated indel signature, ID18 (**Extended Data Fig. 10c-d**). Interestingly, using a  
344 similar methodology, we observed an elevated number of driver mutations assigned to SBS94 and  
345 SBS\_F in Colombia, as well as SBS89 and ID\_J in Argentina, compared to other countries  
346 (**Extended Data Fig. 10e-g**).

347

## 348 **DISCUSSION**

349 Over the last seven decades, colorectal cancer incidence rates have shown complex changes with  
350 marked international variation. Notably, while many high-income countries have seen decreases  
351 in overall incidence rates, there has been an increase amongst adults under the age of 50. If these  
352 trends continue into older age groups, they could reverse the currently overall positive trajectory  
353 for colorectal cancer incidence. In this study, whole-genome sequences of 981 colorectal cancers  
354 from 11 countries revealed evidence of geographic and age-related variation in their landscapes of  
355 somatic mutation, which may contribute to explaining these global trends. These variations were  
356 almost exclusively found in the 802 microsatellite-stable colorectal cancers. For colorectal cancers  
357 with MSI, limited geographic or age-related differences were observed, possibly due to the smaller  
358 sample size and the predominance of somatic mutations resulting from defective DNA repair  
359 mechanisms. Similarly, no differences were noted in colorectal cancers harboring other DNA  
360 repair deficiencies.

361  
362 The prevalence of certain mutational signatures, notably SBS89/DBS8/ID\_J in Argentina and  
363 SBS94/SBS\_F/DBS6 in Colombia, was higher in these countries compared to all others. Although  
364 such geographic variation could, in principle, be due to differences in population-specific  
365 inheritance, it is more plausible that these are due to differences in exogenous environmental or  
366 lifestyle mutagenic exposures. Indeed, apart from country of origin, we also assessed the  
367 variability with genetic ancestry and self-reported ethnicity (**Methods**), although the homogenous  
368 distribution of these characteristics within countries (**Supplementary Fig. 10**) precluded us from  
369 clarifying if the varying prevalence of signatures in different countries was related to genetic or  
370 environmental factors. The natures of the putative exposures underlying SBS89/DBS8/ID\_J and

371 SBS94/SBS\_F/DBS6 are currently unknown. However, SBS89 shares several features with  
372 colibactin-induced signatures SBS88 and ID18. SBS89 has been previously found in normal  
373 colorectal crypts<sup>34</sup> but not in other normal cells. In individuals with SBS89, some crypts have these  
374 mutations while others do not. SBS89 appears to be imprinted on the normal colorectal epithelium  
375 early in life, with mutagenesis ceasing thereafter<sup>34</sup>. Moreover, SBS89 mutations show  
376 transcriptional strand bias<sup>34</sup>, a common trait of mutations caused by exogenous mutagenic  
377 exposures that form bulky covalent DNA adducts. Thus, SBS89 may also be caused by a mutagen  
378 originating from the colorectal microbiome and it is conceivable that multiple microbiome-derived  
379 mutagens may contribute to the mutation burden of the colorectal epithelium. Although the impact  
380 of country-specific microbiome-derived exposures on geographic differences in colorectal cancer  
381 incidence remains unclear, the correlations between colorectal cancer ASR and signatures SBS88  
382 and ID18 suggest that microbiome-derived colibactin exposure may influence colorectal cancer  
383 incidence rates. Nonetheless, further studies are necessary to thoroughly investigate this  
384 hypothesis.

385  
386 The evidence for enrichment of SBS88 and ID18 mutation burdens in early-onset colorectal  
387 cancers may indicate a role for colibactin exposure in the increase in early-onset colorectal cancer  
388 incidence over the last 20 years. Prior studies have indicated that mutagenesis due to colibactin  
389 exposure can occur within the first decade of life and then ceases<sup>34</sup>. In some instances, the mutation  
390 burden caused by this early-life mutation burst can endow affected colorectal crypts with the  
391 equivalent of decades of mutation accumulation and, thus, this ‘head start’ could plausibly result  
392 in an increased risk of early-onset cancers. One mechanism by which colibactin-induced  
393 mutagenesis might contribute to colorectal neoplastic change is by somatically inactivating one



394 copy of *APC* through the generation of protein-truncating driver mutations. Since *APC* mutations  
395 usually occur early in the sequence of driver mutations leading to colorectal cancer<sup>38,44</sup>, a first-hit  
396 inactivating mutation in *APC* during early life could put an individual several decades ahead for  
397 developing colorectal cancer and resulting in a higher likelihood of early-onset colorectal cancer.  
398 The mutation profile of SBS88, with its preponderance of T>C substitutions, is intrinsically  
399 ineffective in generating translation termination codons and SBS88 accounts for only a small  
400 proportion of *APC* driver base substitutions. However, colibactin mutagenesis entails a relatively  
401 high proportion of ID mutations, with the characteristic profile of ID18, almost all of which will  
402 introduce translational frameshifts in coding sequences. ID18 accounts for approximately one  
403 quarter of *APC* indel drivers in colibactin positive cancers and is elevated amongst *APC* indel  
404 drivers compared to indel drivers in other cancer genes such as *TP53*, which occur later in the  
405 multistep process of colorectal carcinogenesis<sup>45</sup>. Thus colibactin-induced indel driver mutations in  
406 *APC* may account for a substantial proportion of any putative impact colibactin exposure has on  
407 colorectal carcinogenesis. Conversely, the unexpected increase in driver mutations observed in  
408 late-onset colorectal cancers might suggest that we failed to identify all driver mutational events  
409 in early-onset cases, possibly overlooking additional effects of colibactin or other mutagenic  
410 exposures, and potentially related to alterations beyond *APC*, as early-onset cases are enriched in  
411 *APC* wild-type tumors<sup>41</sup>. In this context, BMI, diet, lifestyle, and other exposomal factors—  
412 particularly in early life—may play an important mutagenic role, with the lack of analyses on these  
413 factors being a limitation of the current study.

414

415 Although our results show for the first time an association between the presence of colibactin-  
416 induced mutational signatures and early-onset colorectal cancer, complementing the prior finding

417 that tumors harboring colibactin mutagenesis have a younger average age at diagnosis<sup>26</sup>, further  
418 research is required to establish causality. Future studies should examine the SBS88 and ID18  
419 mutation burdens of normal colorectal crypts from individuals with early-onset colorectal cancer  
420 (cases) and age-matched healthy individuals (controls) with the expectation of an enrichment in  
421 cancer cases if colibactin mutagenesis is causally implicated. If so, the increase in early-onset  
422 colorectal cancer over the last 30 years would indicate that an increased exposure to colibactin in  
423 affected populations occurred during the second half of the 20<sup>th</sup> century, perhaps due to increasing  
424 prevalence of *pks*+ bacteria, and genome sequences of appropriately selected colorectal cancers  
425 and normal colorectal tissues would inform on this historical flux. These studies could be  
426 supported by international and, if possible, retrospective studies of the prevalence of colibactin-  
427 producing *pks*+ bacteria in the colorectal microbiome, which should include paired stool samples  
428 or other methods for robust microbiome analysis, not available for the current study. Finally,  
429 definitive evidence of a causal role for colibactin in early-onset colorectal carcinogenesis would  
430 be provided by prevention of early-life exposure to colibactin-producing bacteria reducing cancer  
431 incidence.

432

433 In summary, mutational epidemiology reveals country-specific and age-specific variations in the  
434 prevalence of certain mutational signatures. The results also highlight the potential role of the large  
435 intestine microbiome as an early-life mutagenic factor in the development of colorectal cancer.

436

437 **FIGURE LEGENDS**

438 **Fig. 1. Geographic, clinical, and molecular characterization of the Mutographs colorectal**  
439 **cancer cohort. a**, Geographic distribution of the 981 patients across four continents and 11  
440 countries, with an indication of the total number of cases as well as the percentage of early-onset  
441 (eo) cases below 50 years of age. Countries were colored according to their age-standardized  
442 incidence rates (ASR) per 100,000 individuals. The designations employed and the presentation  
443 of the material in this publication do not imply the expression of any opinion whatsoever on the  
444 part of the authors or their institutions concerning the legal status of any country, territory, city or  
445 area or of its authorities, or concerning the delimitation of its frontiers or boundaries. **b**, Tumor  
446 subsite distribution of the cohort across the colorectum, with an indication of the total number of  
447 cases and the percentage of early-onset cases. Different subsites were colored according to the  
448 percentage of early-onset cases. Additional 2 cases had unspecified subsites. **c**, Scatter plots  
449 indicating the distribution of molecular subgroups across the sequenced tumors according to the  
450 total number of single base substitutions (SBS) and small insertions and deletions (indels; ID), as  
451 well as the percentage of genome aberrated (PGA). Cases for which tumor purity was insufficient  
452 to determine an accurate copy number profile or without large copy number alterations (65/981)  
453 were excluded from the SBS - PGA panel. **d**, Box plots indicating the distribution of SBS and ID  
454 across early-onset (under 50 years; purple) and late-onset (50 years or older; green) microsatellite  
455 stable (MSS) colorectal tumors. Statistically significant differences were evaluated using  
456 multivariable linear regression models adjusted by sex, country, tumor subsite, and tumor purity.  
457 The line within the box is plotted at the median, while the upper and lower ends indicate the 25<sup>th</sup>  
458 and 75<sup>th</sup> percentiles. Whiskers show 1.5 × interquartile range, and values outside it are shown as  
459 individual data points. **e-g**, Average mutational profiles of early-onset and late-onset MSS

460 colorectal tumors for SBS (SBS-288 mutational context; **e**), ID (ID-83 mutational context; **f**), and  
461 copy number alterations (CN-68 mutational context; **g**).

462

463 **Fig. 2. Geographic variation of mutational signatures in microsatellite stable colorectal**

464 **cancers. a**, Dot plot indicating the variation of mutational signature prevalence in specific

465 countries compared to all others. Statistically significant enrichments were evaluated using

466 multivariable logistic regression models adjusted by age of diagnosis, sex, tumor subsite, and

467 tumor purity. Firth's bias-reduced logistic regressions were used for regression presenting

468 complete or quasi-complete separation. Data points were colored according to the odds ratio (OR)

469 of the enrichment, with their size representing statistical significance. P-values were adjusted for

470 multiple comparisons based on the total number of mutational signatures considered per variant

471 type and the total number of countries assessed, and reported as q-values. Q-values<0.05 were

472 considered statistically significant and marked in red. **b-c**, Geographic distribution of the ID\_J (**b**)

473 and SBS\_F (**c**) mutational signatures. Countries were colored based on the signature prevalence.

474 **d**, Volcano plots indicating the association of mutational signature activities with the age-

475 standardized incidence rates. Statistically significant associations were evaluated using

476 multivariable linear regression models adjusted by age of diagnosis, sex, tumor subsite, and tumor

477 purity. P-values were adjusted for multiple comparisons based on the total number of mutational

478 signatures considered per variant type and reported as q-values. Horizontal lines marking

479 statistically significant thresholds were included at 0.05 (dashed orange line) and 0.01 q-values

480 (dashed red line). **e-f**, Scatter plots indicating the association of the mutations attributed to the

481 SBS88 and ID18 mutational signatures with the age-standardized incidence rates (ASR) across

482 countries for colorectal cancer (**e**), and independently for colon and rectal cancers (**f**). Data points

483 were colored based on signature prevalence, with their size indicating the total number of cases  
484 per country. Statistically significant associations were evaluated using the sample-level  
485 multivariable linear regression models used in **d** (**e**), and similar multivariable linear regression  
486 models adjusted by age of diagnosis, sex, and tumor purity (**f**).

487

488 **Fig. 3. Variation of mutational signatures with age of onset in microsatellite stable colorectal**

489 **cancers. a**, Volcano plots indicating the enrichment of mutational signature prevalence in early-

490 onset and late-onset cases. Statistically significant enrichments were evaluated using multivariable

491 logistic regression models for age of onset categorized in two subgroups (early-onset, <50 years

492 of age; and late-onset,  $\geq 50$ ) and adjusted by sex, country, tumor subsite, and tumor purity. Firth's

493 bias-reduced logistic regressions were used for regression presenting complete or quasi-complete

494 separation. P-values were adjusted for multiple comparisons based on the total number of

495 mutational signatures considered per variant type and reported as q-values. Horizontal lines

496 marking statistically significant thresholds were included at 0.05 (dashed orange line) and 0.01 q-

497 values (dashed red line). **b**, Line plots indicating mutational signature prevalence trend across ages

498 of onset, using five different age groups. Signatures significantly enriched in early-onset or late-

499 onset cases (as shown in **a**) were colored in purple and green, respective, whereas signatures not

500 varying significantly with age were colored in grey. **c**, Bar plots indicating mutational signature

501 prevalence across age groups, with indication of the total number of cases where signatures were

502 detected. Statistically significant trends were evaluated using multivariable logistic regression

503 models for age categorized in five subgroups (0-39, 40-49, 50-59, 60-69,  $\geq 70$ ) and adjusted by

504 sex, country, tumor subsite, and tumor purity. Firth's bias-reduced logistic regressions were used

505 for regressions presenting complete or quasi-complete separation. **d-e**, Box plots indicating the

506 variation in age of onset according to the presence of colibactin mutational signatures (either  
507 SBS88, ID18, or both) in all microsatellite stable cases (**d**) and across tumor subsites (**e**).  
508 Statistically significant differences were evaluated using multivariable linear regression models  
509 adjusted by sex, country, tumor purity, and tumor subsite (only for the analysis of all cases, **d**).  
510 The line within the box is plotted at the median, while the upper and lower ends indicate the 25<sup>th</sup>  
511 and 75<sup>th</sup> percentiles. Whiskers show  $1.5 \times$  interquartile range, and values outside it are shown as  
512 individual data points.

513

514 **Fig. 4. Colibactin mutagenesis as an early event in microsatellite stable colorectal cancer**

515 **evolution. a**, Box plots indicating the fold-change of the relative contribution per sample of each  
516 signature between early clonal and late clonal single base substitutions (SBS, left) and small  
517 insertions and deletions (ID, right). SBS signatures that generated early clonal SBSs in fewer than  
518 50 samples and also generated late clonal SBSs in fewer than 50 samples were excluded from the  
519 analysis. Similarly, ID signatures that generated early clonal IDs in fewer than 20 samples and late  
520 clonal IDs in fewer than 20 samples were also excluded. Signatures were sorted by their median  
521 fold-change. The line within the box is plotted at the median, while the upper and lower ends  
522 indicate the 25<sup>th</sup> and 75<sup>th</sup> percentiles. Whiskers show  $1.5 \times$  interquartile range, and values outside  
523 it are shown as individual data points. **b**, Bar plot indicating the lack of concordance between  
524 colibactin exposure status determined by the presence of colibactin-induced mutational signatures  
525 SBS88 or ID18, and the microbiome *pks* status. Statistical significance was evaluated using a  
526 multivariable Firth's bias-reduced logistic regression model (due to quasi-complete separation)  
527 adjusted by age of diagnosis, sex, country, tumor subsite, and tumor purity. **c-d**, Distribution of  
528 the age of onset (**c**) and cases across age groups (**d**) based on the detection of colibactin-positive

529 samples using genomic and microbiome status. The genomic status is defined by the presence of  
530 SBS88 or ID18 signatures, while the microbiome status (*pks*) is determined by coverage of at least  
531 half of the *pks* island, and suggests ongoing or active *pks*<sup>+</sup> bacterial infection. Statistical  
532 significance was evaluated using a multivariable linear regression model adjusted by sex, country,  
533 tumor subsite, and tumor purity.

534

535 **Fig. 5. Variation of driver mutations with age of onset and association with colibactin**

536 **mutagenesis in microsatellite stable colorectal cancers. a**, Bar plot indicating the prevalence of

537 driver mutations affecting the 48 bioinformatically detected driver genes in microsatellite stable

538 colorectal cancers. Genes were colored according to their status as known cancer driver genes for

539 colorectal cancer, known cancer driver genes for other cancer types, or newly detected cancer

540 driver genes. **b**, Box plots indicating the distribution of total driver mutations across early-onset

541 (under 50 years of age; purple) and late-onset (50 or over; green) tumors. Statistical significance

542 was evaluated using a multivariable linear regression model adjusted by sex, country, tumor

543 subsite, and tumor purity. The line within the box is plotted at the median, while the upper and

544 lower ends indicate the 25<sup>th</sup> and 75<sup>th</sup> percentiles. Whiskers show 1.5 × interquartile range, and

545 values outside it are shown as individual data points. **c**, Volcano plot indicating the enrichment of

546 driver mutations in cancer driver genes in early-onset and late-onset cases. Statistically significant

547 enrichments were evaluated using multivariable logistic regression models adjusted by sex,

548 country, tumor subsite, and tumor purity. Firth's bias-reduced logistic regressions were used for

549 regressions presenting complete or quasi-complete separation. P-values were adjusted for multiple

550 comparisons based on the total number of cancer driver genes considered and reported as q-values.

551 Horizontal lines marking statistically significant thresholds were included at 0.05 (dashed orange

552 line) and 0.01 q-values (dashed red line). **d**, Line plot indicating the prevalence of driver mutations  
553 in cancer driver genes across ages of onset, using five different age groups. Cancer driver genes  
554 significantly enriched in late-onset cases (as shown in **c**) were colored in green, whereas genes not  
555 varying significantly with age of onset were colored in grey. **e**, Bar plots indicating the proportion  
556 and number of driver mutations probabilistically assigned to colibactin-induced and other  
557 mutational signatures, including single base substitutions (**e**) and small insertions and deletions  
558 (indels; **f**). Driver mutations were divided into different groups, including the *APC* c.835-8A>G  
559 splicing-associated driver mutation, other *APC* driver mutations, *TP53* driver mutations, and driver  
560 mutations affecting other cancer driver genes.



561 **EXTENDED DATA FIGURE LEGENDS**

562 **Extended Data Fig. 1. Mutational profiles across molecular subtypes and ages of onset. a-b,**

563 Average mutational profiles of microsatellite stable (MSS; **a**) and microsatellite unstable (MSI; **b**)  
564 colorectal tumors for single base substitutions (SBS-288 mutational context), small insertions and  
565 deletions (ID-83 mutational context), doublet base substitutions (DBS-78 mutational context),  
566 copy number alterations (CN-68 mutational context), and structural variants (SV-38 mutational  
567 context). **c-d**, Average mutational profiles of early-onset and late-onset MSS colorectal tumors for  
568 doublet base substitutions (**c**) and structural variants (**d**).

569

570 **Extended Data Fig. 2. Geographic distribution of mutation burden.** Box plots indicating the

571 distribution of single base substitutions (SBS), small insertions and deletions (ID), doublet base  
572 substitutions (DBS), copy number alterations (CN), and structural variants (SV) across countries  
573 for microsatellite stable (MSS) colorectal tumors. Box plots and data points representing total  
574 number of mutations for each variant type were colored according to each country's colorectal  
575 cancer age-standardized incidence rates (ASR) per 100,000 individuals. A horizontal blue line  
576 indicates the median mutation burden for each variant type. Statistically significant differences  
577 were evaluated using multivariable linear regression models comparing each country to all others  
578 and adjusted by age of diagnosis, sex, tumor subsite, and tumor purity. P-values were adjusted for  
579 multiple comparisons based on the total number of countries assessed and reported as q-values.  
580 The line within the box is plotted at the median, while the upper and lower ends indicate the 25<sup>th</sup>  
581 and 75<sup>th</sup> percentiles. Whiskers show  $1.5 \times$  interquartile range, and values outside it are shown as  
582 individual data points.

583

584 **Extended Data Fig. 3. Geographic distribution of mutational profiles. a-e**, Average mutational  
585 profiles of microsatellite stable (MSS) colorectal tumors for single base substitutions (SBS-288  
586 mutational context; **a**), small insertions and deletions (ID-83 mutational context; **b**), doublet base  
587 substitutions (DBS-78 mutational context; **c**), copy number alterations (CN-68 mutational context;  
588 **d**), and structural variants (SV-38 mutational context; **e**).

589  
590 **Extended Data Fig. 4. Mutational signatures of small mutational events identified in**  
591 **microsatellite stable colorectal cancers. a-c**, Mutational profiles of single base substitution  
592 (SBS) signatures, including COSMICv3.4 reference signatures (**a**), previously reported signatures  
593 not present in COSMIC (**b**), and novel signature SBS\_O (**c**). **d-e**, Mutational profiles of small  
594 insertions and deletions (ID) signatures, including COSMICv3.4 signatures (**d**) and novel  
595 signature ID\_J (**e**). **f**, Mutational profiles of doublet base substitution (DBS) signatures, all  
596 previously reported in COSMIC.

597  
598 **Extended Data Fig. 5. Mutational signatures of large mutational events identified in**  
599 **microsatellite stable colorectal cancers. a-b**, Mutational profiles of copy number (CN)  
600 signatures, including COSMICv3.4 reference signatures (**a**) and novel signature CN\_F (**b**). **c-d**,  
601 Mutational profiles of structural variant (SV) signatures, including COSMIC signatures (**c**) and  
602 novel signatures SV\_B and SV\_D (**d**).

603  
604 **Extended Data Fig. 6. Association of mutational signatures with colorectal, colon, and rectal**  
605 **cancer incidence rates. a-b**, Scatter plots indicating the association of the mutations attributed to  
606 signatures SBS1, SBS\_H, and CN\_F with the age-standardized incidence rates across countries

607 for colorectal cancers (**a**), and independently for colon and rectal cancer (**b**). Data points were  
608 colored based on signature prevalence, with their size indicating the total number of cases per  
609 country. Statistically significant associations were evaluated using the sample-level multivariable  
610 linear regression models used in **Fig. 2d (a)**, and similar multivariable linear regression models  
611 adjusted by age of diagnosis, sex, and tumor purity (**b**). **c**, Bar plots indicating mutational signature  
612 prevalence enrichment between low and high ASR countries (defined as those below or above an  
613 ASR of 7 per 100,000 people, for early-onset colorectal cancer, diagnosed between 20 and 49 years  
614 old). Statistically significant associations were evaluated using multivariable logistic regression  
615 models for early-onset colorectal cancer ASR adjusted by age of diagnosis, sex, tumor subsite, and  
616 tumor purity.

617

618 **Extended Data Fig. 7. Enrichment of colibactin mutagenesis in early-onset colorectal cancers**

619 **based on motif analysis. a**, Box plots indicating the percentage of total W[T>N]W mutations with  
620 the WAWW[T>N]W motif across different age groups. Statistically significant trend was  
621 evaluated using a multivariable linear regression model adjusted by sex, country, tumor subsite,  
622 and tumor purity. The line within the box is plotted at the median, while the upper and lower ends  
623 indicate the 25<sup>th</sup> and 75<sup>th</sup> percentiles. Whiskers show  $1.5 \times$  interquartile range, and values outside  
624 it are shown as individual data points. **b**, Box plots indicating the percentage of total W[T>N]W  
625 mutations with the WAWW[T>N]W motif across samples grouped by colibactin exposure status,  
626 determined by the presence of signatures SBS88 or ID18. Statistical significance was evaluated  
627 using a multivariable linear regression model adjusted by age, sex, country, tumor subsite, and  
628 tumor purity. **c**, Bar plots indicating the prevalence of colibactin exposure across age groups, with  
629 indication of the total number of cases where colibactin signatures were detected. Statistically

630 significant trend was evaluated using a multivariable logistic regression model adjusted by sex,  
631 country, tumor subsite, and tumor purity.

632

633 **Extended Data Fig. 8. Enrichment of colibactin mutagenesis as an early clonal event in early-**

634 **onset and late-onset colorectal cancers. a,** Box plots indicating the fold-change of the relative

635 contribution per sample of each signature between clonal and subclonal single base substitutions

636 (SBS, left) and small insertions and deletions (ID, right). Signatures that generated clonal somatic

637 mutations in fewer than 10 samples and also generated subclonal somatic mutations in fewer than

638 10 samples were excluded from the analysis. The line within the box is plotted at the median, while

639 the upper and lower ends indicate the 25<sup>th</sup> and 75<sup>th</sup> percentiles. Whiskers show 1.5 × interquartile

640 range, and values outside it are shown as individual data points. **b,** Boxplots indicating the fold-

641 change of the relative contribution per sample of each signature between early clonal and late

642 clonal SBS (left) and ID (right) with samples separated by age of diagnosis in early-onset (under

643 50 years of age; purple) and late-onset (50 or over; green). As in **Fig. 4a**, SBS signatures that

644 generated early clonal SBSs in fewer than 50 samples and late clonal SBSs in fewer than 50

645 samples, as well as ID signatures generating early clonal IDs in fewer than 20 samples and late

646 clonal IDs in fewer than 20 samples, were excluded from the analysis.

647

648 **Extended Data Fig. 9. Representative microbiome and genomic profiles of colibactin-exposed**

649 **samples. a-d,** Microbiome and genomic profiles of representative samples corresponding to the

650 four different sample types according to colibactin exposure: genomic<sup>+</sup> and *pks*<sup>+</sup> (**a**), genomic<sup>+</sup>

651 and *pks*<sup>-</sup> (**b**), genomic<sup>-</sup> and *pks*<sup>+</sup> (**c**), and genomic<sup>-</sup> and *pks*<sup>-</sup>. The genomic status is defined by the

652 presence of SBS88 or ID18 signatures, while the microbiome status (*pks*) is determined by the

653 coverage of at least half of the *pks* island, and suggests ongoing and/or active *pks*+ bacterial  
654 infection. Circos plots display Reads Per Kilobase of transcript per Million (RPKM) values across  
655 *clb* genes within the *pks* island (left). Bar plots represent the proportion of mutations attributed to  
656 SBS88 and ID18 colibactin signatures compared to others (center), and are displayed next to  
657 mutational profiles of single base substitutions (SBS-288 mutational context) and small insertions  
658 and deletions (ID-83 mutational context) for each sample (right).

659

660 **Extended Data Fig. 10. Driver mutations associated with colibactin mutagenesis in early-**  
661 **onset and late-onset colibactin positive cases and with country-enriched mutational**  
662 **signatures in microsatellite stable colorectal cancers. a-b,** Bar plots indicating the proportion  
663 and number of driver mutations probabilistically assigned to colibactin-induced and other  
664 mutational signatures, including single base substitutions (**a**) and small insertions and deletions  
665 (indels; **b**), with samples separated by age of diagnosis in early-onset (under 50 years of age; left)  
666 and late-onset (50 or over; right). Driver mutations were divided into different groups, including  
667 the *APC* c.835-8A>G splicing-associated driver mutation, other *APC* driver mutations, *TP53*  
668 driver mutations, and driver mutations affecting other cancer driver genes. **c-d,** Bar plots indicating  
669 the proportion and number of mutations in chromatin modifier genes probabilistically assigned to  
670 colibactin-induced and other mutational signatures, including single base substitutions (**c**) and  
671 indels (**d**), in the 169 colibactin positive cases. **e-g,** Bar plots indicating the proportion and number  
672 of driver single base substitutions (**e** and **f**) and indels (**g**) in cancer driver genes probabilistically  
673 assigned to specific mutational signatures in cases from Colombia (**e**) or Argentina (**f** and **g**)  
674 compared to other countries.

675 **SUPPLEMENTARY FIGURE LEGENDS**

676 **Supplementary Fig. 1. Mutational profiles of homologous recombination deficient colorectal**  
677 **cancers. a**, Mutational profiles of individual samples identified as homologous recombination  
678 deficient for single base substitutions (SBS-288 mutational context) and small insertions and  
679 deletions (ID-83 mutational context). **b**, Mutational signatures previously associated with  
680 homologous recombination deficiency in COSMICv3.4 (SBS3 and ID6).

681  
682 **Supplementary Fig. 2. Mutational profiles of base excision repair deficient and *POLD1***  
683 **mutated colorectal cancers. a**, Mutational profile of an individual sample identified as base  
684 excision repair deficient due to mutations in *MUTYH* for single base substitutions (SBS-288  
685 mutational context). **b**, Mutational signature previously associated with base excision repair  
686 deficiency due to mutations in *MUTYH* in COSMICv3.4 (SBS36). **c**, Mutational profiles of  
687 individual samples identified as base excision repair deficient due to mutations in *NTHL1* for  
688 single base substitutions (SBS-288 mutational context). **d**, Mutational signature previously  
689 associated with base excision repair deficiency due to mutations in *NTHL1* in COSMICv3.4  
690 (SBS30). **e**, Mutational profile of an individual sample identified as base excision repair deficient  
691 due to mutations in *OGGI* for single base substitutions (SBS-288 mutational context). **f**,  
692 Mutational signature previously associated with base excision repair deficiency due to mutations  
693 in *OGGI* (Signal signature SBS108). **g**, Mutational profiles of individual samples harboring  
694 mutations in *POLD1* for single base substitutions (SBS-288 mutational context). **h**, Mutational  
695 signature previously associated with mutations in *POLD1* in COSMICv3.4 (SBS10c).

696

697 **Supplementary Fig. 3. Mutational profiles of *POLE* mutated colorectal cancers.** **a**, Mutational  
698 profiles of individual samples harboring mutations in *POLE* for single base substitutions (SBS-  
699 288 mutational context). **b**, Mutational signatures previously associated with mutations in *POLE*  
700 in COSMICv3.4 (SBS10a, SBS10b, and SBS28).

701  
702 **Supplementary Fig. 4. Mutational profiles of cases treated with chemotherapy for prior**  
703 **cancers.** **a**, Mutational profiles of individual samples treated with chemotherapy for prior cancers  
704 for single base substitutions (SBS-288 mutational context) and doublet base substitutions (DBS-  
705 78 mutational context). **b**, Mutational signatures previously associated with chemotherapy in  
706 COSMICv3.4 (SBS25, SBS31, SBS35, and DBS5).

707  
708 **Supplementary Fig. 5. Prevalence of mutational signatures in microsatellite stable colorectal**  
709 **cancers by country.** **a-e**, SBS (**a**), ID (**b**), DBS (**c**), CN (**d**), and SV (**e**) signatures.

710  
711 **Supplementary Fig. 6. Main age association analyses adjusted by self-reported ethnicity**  
712 **instead of country of origin.** **a-e**, Replicates of Fig. 1d (**a**), Fig.3a (**b**), and Fig. 3c-e (**c-e**).

713  
714 **Supplementary Fig. 7. Main age association analyses adjusted by the first five principal**  
715 **components of the genetic ancestry analysis instead of country of origin.** **a-e**, Replicates of  
716 Fig. 1d (**a**), Fig.3a (**b**), and Fig. 3c-e (**c-e**).

717

718 **Supplementary Fig. 8. Main age association analyses adjusted by genetic ancestry groups**  
719 **(ADMIX, AFR, EAS, EUR) instead of country of origin. a-e**, Replicates of Fig. 1d (a), Fig.3a  
720 (b), and Fig. 3c-e (c-e).

721  
722 **Supplementary Fig. 9. Variation of mutational signatures associated with earlier age of onset**  
723 **with tumor subsite in microsatellite stable colorectal cancers.**

724  
725 **Supplementary Fig. 10. Genetic ancestry (a) and self-reported ethnicity (b) distribution by**  
726 **country.**

727  
728 **Supplementary Fig. 11. Comparison of clinicopathological characteristics at baseline by**  
729 **country, comparing Mutographs samples vs. GLOBOCAN-based expectations.**

730  
731 **Supplementary Fig. 12. Mutational signature reconstruction of an individual microsatellite**  
732 **unstable tumor not validated by droplet digital PCR.** Multiple mutational signatures were  
733 assigned to this case, as indicated in **Supplementary Note Table 8**, including microsatellite  
734 instability-associated COSMICv3.4 signature SBS15 and SBS26, suggesting the presence of  
735 microsatellite instability.

736  
737 **Supplementary Fig. 13. Novel mutational signatures identified in microsatellite unstable**  
738 **colorectal cancers. a-b**, Mutational profiles of single base substitution (SBS) signatures not  
739 matching any previous COSMICv3.4 signature, including three novel signatures (a) and a  
740 previously reported signature (b). c, Mutational profile of a novel doublet base substitution (DBS)



741 signature not previously reported in COSMIC. **d-e**, Exemplar mutational profiles of individual  
742 colorectal cancers supporting the three novel signatures indicated in **a (d)** and the previously  
743 reported signature in **b (e)**.

744  
745 **Supplementary Fig. 14. Variation of germline pathogenic variants with age of onset. a-b,**  
746 Volcano plots indicating the enrichment of pathogenic or likely pathogenic germline variants in  
747 early-onset and late-onset cases in microsatellite stable (MSS; **a**) and unstable colorectal cancers  
748 (MSI; **b**). Separate analyses were performed for all individual genes (left), and for genes grouped  
749 in colorectal cancer predisposition syndromes (Lynch syndrome and Cowden syndrome),  
750 homologous recombination, or DNA damage repair-associated genes (right). Statistically  
751 significant enrichments were evaluated using multivariable logistic regression models for age of  
752 onset categorized in two subgroups (early-onset, <50 years of age; and late-onset,  $\geq 50$ ) and  
753 adjusted by sex, country, and tumor subsite, and tumor purity. Firth's bias-reduced logistic  
754 regressions were used for regressions presenting complete or quasi-complete separation. P-values  
755 were adjusted for multiple comparisons based on the total number of germline variants considered  
756 and reported as q-values. Horizontal lines marking statistically significant thresholds were  
757 included at 0.05 (dashed orange line) and 0.01 q-values (dashed red line).

758  
759 **Supplementary Fig. 15. Dendrogram of the relative abundances of the considered bacterial**  
760 **genus based on the Bray-Curtis distance.**

761  
762 **Supplementary Fig. 16. Single base substitution mutational signatures extracted by**  
763 **SigProfilerExtractor using the SBS-288 and SBS-1536 mutational contexts in microsatellite**

764 **stable colorectal cancers.** All single base substitution (SBS) *de novo* signatures extracted using  
765 the SBS-288 (16 signatures) and SBS-1536 (14 signatures) mutational contexts, shown side by  
766 side for comparison. Equivalent signatures were not extracted in SBS-1536 format for SBS288E  
767 and SBS288K. For clarity, the signature context is retained in the signature names in this figure.  
768 The extended context for SBS-1536 signatures is omitted from the figure. Instead, the SBS-96  
769 down-sampled version of the SBS-1536 *de novo* extracted signatures was used to display the  
770 signatures.

771  
772 **Supplementary Fig. 17. Single base substitution mutational signatures extracted by**  
773 **mSigHdp in microsatellite stable colorectal cancers.** Seventeen single base substitution (SBS)  
774 *de novo* signatures were extracted by mSigHdp, using the SBS-96 mutational context.

775  
776 **Supplementary Fig. 18. Small insertion and deletion mutational signatures extracted by**  
777 **SigProfilerExtractor and mSigHdp in microsatellite stable colorectal cancers. a,** Ten small  
778 insertions and deletions (ID) *de novo* signatures were extracted by SigProfilerExtractor using the  
779 ID-83 mutational context. **b,** Eight ID *de novo* signatures were extracted by mSigHdp.

780  
781 **Supplementary Fig. 19. Sensitivity analysis for the detection of colibactin signatures in**  
782 **microsatellite stable colorectal cancers. a-b,** Sensitivity analysis for SBS88 (a) and ID18 (b)  
783 detection, including average relative activity of the signature detected across all colibactin negative  
784 samples and simulations (top), activity of the signature per sample across all colibactin negative  
785 samples and simulations (middle), and activity of the signature per sample across all colibactin  
786 positive samples compared to the median detection of the signature in the simulated data (bottom).

787 **SUPPLEMENTARY TABLES**

788 **Supplementary Table 1. Summary of incidence rates, sex, age, tumor subsite, and molecular**  
789 **subgroups across countries included in the Mutographs colorectal cancer cohort.**

790

791 **Supplementary Table 2. Germline mutations in mismatch repair genes in MSI Lynch**  
792 **syndrome cases.**

793

794 **Supplementary Table 3. Germline and somatic mutations in DNA repair deficient cases.**

795

796 **Supplementary Table 4. Detection results of homologous recombination deficiency with**  
797 **CHORD.**

798

799 **Supplementary Table 5. Mutational profiles of *de novo* SBS signatures extracted in MSS**  
800 **colorectal cancer cases.**

801

802 **Supplementary Table 6. Mutational profiles of *de novo* ID signatures extracted in MSS**  
803 **colorectal cancer cases.**

804

805 **Supplementary Table 7. Mutational profiles of *de novo* DBS signatures extracted in MSS**  
806 **colorectal cancer cases.**

807

808 **Supplementary Table 8. Mutational profiles of *de novo* CN signatures extracted in MSS**  
809 **colorectal cancer cases.**

810

811 **Supplementary Table 9. Mutational profiles of *de novo* SV signatures extracted in MSS**  
812 **colorectal cancer cases.**

813

814 **Supplementary Table 10. Decomposition of *de novo* MSS colorectal cancer signatures into**  
815 **previously reported signatures.**

816

817 **Supplementary Table 11. Sample attributions of decomposed SBS signatures in MSS**  
818 **colorectal cancers.**

819

820 **Supplementary Table 12. Sample attributions of decomposed ID signatures in MSS**  
821 **colorectal cancers.**

822

823 **Supplementary Table 13. Sample attributions of decomposed DBS signatures in MSS**  
824 **colorectal cancers.**

825

826 **Supplementary Table 14. Sample attributions of decomposed CN signatures in MSS**  
827 **colorectal cancers.**

828

829 **Supplementary Table 15. Sample attributions of decomposed SBS signatures in SV**  
830 **colorectal cancers.**

831

832 **Supplementary Table 16. Enrichment of mutational signature prevalence in specific**  
833 **countries compared to all others in MSS colorectal cancers.**

834

835 **Supplementary Table 17. Enrichment of mutational signature activities with colorectal**  
836 **cancer incidence in MSS colorectal cancers.**

837

838 **Supplementary Table 18. Enrichment of mutational signature activities with colon and rectal**  
839 **cancer incidence in MSS colon and rectal cancers.**

840

841 **Supplementary Table 19. Enrichment of mutational signature prevalence in late-onset**  
842 **compared to early-onset MSS colorectal cancers.**

843

844 **Supplementary Table 20. Trend enrichment of mutational signature prevalence with age of**  
845 **diagnosis in MSS colorectal cancers.**

846

847 **Supplementary Table 21. Driver genes detected in MSS colorectal cancers.**

848

849 **Supplementary Table 22. Enrichment of driver mutations in cancer driver genes in late-**  
850 **onset compared to early-onset MSS colorectal cancers.**

851

852 **Supplementary Table 23. Enrichment of hotspot driver mutations in late-onset compared to**  
853 **early-onset MSS colorectal cancers.**

854

855 **Supplementary Table 24. Enrichment of driver mutations in cancer driver genes in specific**  
856 **countries compared to all others in MSS colorectal cancers.**

857

858 **Supplementary Table 25. Enrichment of hotspot driver mutations in specific countries**  
859 **compared to all others in MSS colorectal cancers.**

860

861 **Supplementary Table 26. Details of individual case collections.**

862

863 **Supplementary Table 27. Clinicopathological characteristics of included and excluded cases.**

## 864 **ONLINE METHODS**

### 865 **Recruitment of patients and informed consent**

866 The International Agency for Research on Cancer (IARC/WHO) coordinated case recruitment  
867 through an international network of 17 collaborators from 11 participating countries in North  
868 America, South America, Asia, and Europe (**Supplementary Table 26**). The inclusion criteria for  
869 patients were  $\geq 18$  years of age (ranging from 18 to 95, with a mean of 64 and a standard deviation  
870 of 12), confirmed diagnosis of primary colorectal cancer, and no prior treatment for colorectal  
871 cancer. Informed consent was obtained for all participants. Patients were excluded if they had any  
872 condition that could interfere with their ability to provide informed consent or if there were no  
873 means of obtaining adequate tissues or associated data as per the protocol requirements. Ethical  
874 approvals were first obtained from each Local Research Ethics Committee and Federal Ethics  
875 Committee when applicable, as well as from the IARC/WHO Ethics Committee.

876

### 877 **Bio-samples and data collection**

878 Dedicated standard operating procedures, following guidelines from the International Cancer  
879 Genome Consortium (ICGC), were designed by IARC/WHO to select appropriate case series with  
880 complete biological samples and exposure information<sup>46</sup>, as described previously<sup>5,6,8</sup>  
881 (**Supplementary Table 26**). In brief, for all case series included, anthropometric measures were  
882 taken, together with relevant information regarding medical and familial history. All biological  
883 samples from retrospective cohorts were collected using rigorous, standardized protocols and  
884 fulfilled the required standards of sample collection defined by the IARC/WHO for sequencing  
885 and analysis. Potential limitations of using retrospective clinical data collected using different  
886 protocols from different populations were addressed by a central data harmonization to ensure a

887 comparable group of exposure variables (**Supplementary Table 26**). All patient-related data were  
888 pseudonymized locally through the use of a dedicated alpha-numerical identifier system before  
889 being transferred to the IARC/WHO central database.

890

### 891 **Expert pathology review**

892 Original diagnostic pathology departments provided diagnostic histological details of contributing  
893 cases through standard abstract forms, together with a representative hematoxylin-eosin-stained  
894 slide of formalin-fixed paraffin-embedded tumor tissues whenever possible. IARC/WHO  
895 centralized the entire pathology workflow and coordinated a centralized digital pathology  
896 examination of the frozen tumor tissues collected for the study as well as formalin-fixed paraffin-  
897 embedded sections when available, via a web-based approach and dedicated expert panel  
898 following standardized procedures as described previously<sup>5,6</sup>. A minimum of 50% viable tumor  
899 cells was required for eligibility for whole-genome sequencing. In summary, frozen tumor tissues  
900 were first examined to confirm the morphological type and the percentage of viable tumor cells.  
901 A random selection of tumor tissues was independently evaluated by a second pathologist.  
902 Enrichment of tumor component was performed by dissection of the non-tumoral part, if  
903 necessary.

904

### 905 **DNA extraction**

906 A total of 1,977 primary colorectal cancer patients were enrolled into the study, including  
907 biological samples for 1,946 cases and sequencing data (FASTQ) for 31 cases from Japan. Of  
908 these, 906 samples (45.8%) were excluded due to insufficient viable tumor cells (pathology level)  
909 or inadequate DNA (tumor or germline). Extraction of DNA from fresh frozen primary tumor and



910 matched blood/normal tissue samples was centrally conducted at IARC/WHO (except for samples  
911 from Japan) following a similarly standardized DNA extraction procedure. Germline DNA was  
912 extracted from whole blood ( $n=1,015$ ), except for a small subset of Canadian cases ( $n=25$ ) where  
913 only adjacent normal tissue was available, following previously described protocols and  
914 methods<sup>5,6</sup>. As a result, DNA from 1,040 cases was sent to the Wellcome Sanger Institute for  
915 whole-genome sequencing.

916

### 917 **Whole-genome sequencing**

918 Fluidigm SNP genotyping with a custom panel was performed to ensure that each pair of tumor  
919 and matched normal samples originated from the same individual. Whole-genome sequencing  
920 (150bp paired-end) was performed on the Illumina NovaSeq 6000 platform with a target coverage  
921 of 40x for tumors and 20x for matched normal tissues. All sequencing reads were aligned to the  
922 GRCh38 human reference genome using the Burrows-Wheeler Aligner MEM (BWA-MEM;  
923 v0.7.16a and v0.7.17)<sup>47</sup>. Post-sequencing quality control metrics were applied for total coverage,  
924 evenness of coverage, contamination, and total number of somatic single base substitutions  
925 (SBSs). Cases were excluded if coverage was below 30x for tumor or 15x for normal tissue. For  
926 evenness of coverage, the median over mean coverage (MoM) score was calculated. Tumors with  
927 MoM scores outside the range of values determined by previous studies<sup>48</sup> to be appropriate for  
928 whole-genome sequencing (0.92-1.09) were excluded. Conpair<sup>49</sup> was used to detect  
929 contamination, cases were excluded if the result was greater than 3%<sup>48</sup>. Lastly, samples with  
930 <1,000 total somatic SBSs were also excluded. A total of 981 pairs of colorectal cancer and  
931 matched-normal tissue passed all criteria. Comparing the clinicopathological characteristics  
932 between the included and excluded patients revealed very similar traits (**Supplementary Table**

933 27), and comparable to those expected for each country according to GLOBOCAN metrics  
934 (obtained from <https://gco.iarc.who.int/today/en/dataviz/>; **Supplementary Fig. 11**).

935

### 936 **Germline variant calling**

937 Germline SNVs and indels were derived from whole-genome sequencing from the normal paired  
938 material for each individual using Strelka2 with appropriate quality-control criteria<sup>50</sup>. Variant calls  
939 were then derived into genotypes for each individual and annotated using ANNOVAR<sup>51</sup>.

940

### 941 **Somatic variant calling**

942 Variant calling was performed using the standard Sanger bioinformatics analysis pipeline  
943 (<https://github.com/cancerit>). Copy number profiles were determined using ASCAT<sup>52</sup> and  
944 BATTENBERG<sup>53</sup> when tumor purity allowed. SNVs were called with cgpCaVEMan<sup>54</sup>, indels  
945 were called with cgpPINDEL<sup>55</sup>, and structural rearrangements were called using BRASS  
946 (<https://github.com/cancerit/BRASS>). CaVEMan and BRASS were run using the copy number  
947 profile and purity values determined from ASCAT when possible (complete pipeline,  $n=916$ ).  
948 When tumor purity was insufficient to determine an accurate copy number profile (partial pipeline,  
949  $n=31$ ) CaVEMan and BRASS were run using copy number defaults and an estimate of purity  
950 obtained from ASCAT. Finally, for a subset of cases which had no large copy number alterations  
951 (copy number normal pipeline,  $n=34$ ), CaVEMan and BRASS were run using copy number  
952 defaults and an estimate of purity calculated by the median variant allele frequency (VAF) of indels  
953 multiplied by two. For SNVs, additional filters ( $ASRD \geq 140$  and  $CLPM=0$ ) were applied in  
954 addition to the standard PASS filter to remove potential false positive calls. To further exclude the  
955 possibility of caller-specific artifacts being included in the analysis, a second variant caller,

956 Strelka2<sup>50</sup>, was run for SNVs and indels. Only variants called by both the Sanger variant calling  
957 pipeline and Strelka2 were included in subsequent analysis.

958

### 959 **Generation of mutational matrices**

960 Mutational matrices for single base substitutions (SBS), indels (ID), doublet base substitutions  
961 (DBS), copy number alterations (CN), and structural variants (SV) were generated using  
962 SigProfilerMatrixGenerator with default options (v1.2.0)<sup>56,57</sup>.

963

### 964 **Microsatellite instability validation**

965 The presence of microsatellite instability (MSI) in colorectal cancers was validated using the  
966 QX200 Droplet Digital PCR System (Bio-Rad, Hercules, CA, USA) for the detection of five  
967 microsatellite markers (BAT25, BAT26, NR21, NR24, and Mono27) commercially pooled in three  
968 primer–probe mix assays, as previously described<sup>58</sup>. Briefly, samples were tested in duplicate, and  
969 each reaction comprised 1× ddPCR Multiplex Supermix for probes (Bio-Rad), 1X primer–probe  
970 mix, and 10 ng of extracted tumor DNA, in a total volume of 22 µl. MSI-positive, negative, and  
971 no-template (nuclease-free water) controls were included in each experiment. Droplet generation  
972 and plate preparation for thermal cycling amplification were performed using the QX200 AutoDG  
973 Droplet Digital PCR System (Bio-Rad). The following PCR protocol was applied on a C1000  
974 Touch Thermal Cycler (Bio-Rad): 37 °C for 30 min, 95 °C for 10 min, followed by 40 cycles of  
975 denaturation at 94 °C for 30 s, annealing at 55 °C for 1 min, with a final extension at 98 °C for 10  
976 min. Following PCR amplification, fluorescence signals were quantified using the QX200 Droplet  
977 Reader (Bio-Rad), and data were analyzed with QuantaSoft Analysis Pro v.1.0.596.0525 (Bio-  
978 Rad) software. Positive and negative controls served as guides to call markers and delineate

979 clusters. For each assay, the cluster at the bottom left of the x–y plot was designated as the negative  
980 population. Clusters located vertically and horizontally from the negative cluster were identified  
981 as the mutant population, while clusters located diagonally from the negative cluster represented  
982 the wild-type population. Tumors were characterized for the MSI phenotype by analyzing the  
983 results for all five markers using the following criteria: MSI positive if two or more mutant  
984 microsatellite markers were observed, and microsatellite stable (MSS; *i.e.*, MSI negative) when  
985 none or only one of the microsatellite markers was altered (**Supplementary Note**).

986

### 987 **Extraction and decomposition of mutational signatures**

988 Mutational signatures were primarily extracted using SigProfilerExtractor<sup>35</sup>, based on nonnegative  
989 matrix factorization, and validated by mSigHdp<sup>59</sup>, based on hierarchical Dirichlet process mixture  
990 models. For SigProfilerExtractor (v1.1.21), *de novo* mutational signatures were extracted from  
991 SBS, DBS, and ID mutational matrices using 500 NMF replicates (nmf\_replicates=500),  
992 nndsvd\_min initialization (nmf\_init="nndsvd\_min"), and default parameters. Extractions were  
993 performed separately on the subsets of 802 MSS and 153 MSI cases (**Supplementary Note**). *De*  
994 *novo* SBS mutational signatures were extracted for both SBS-288 and SBS-1536 contexts, which,  
995 beyond the common SBS-96 trinucleotide context using the mutated base and the 5' and 3'  
996 adjacent nucleotides<sup>57,60</sup>, also consider the transcriptional strand bias and the pentanucleotide  
997 context (two 5' and 3' adjacent nucleotides), respectively<sup>57</sup>. The results were largely concordant,  
998 with the SBS-288 *de novo* signatures allowing additional separation of mutational processes. The  
999 SBS-1536 results can be found in the **Supplementary Note**. Therefore, the SBS-288 *de novo*  
1000 signatures were taken forward for further analysis (**Supplementary Table 5**). Previously

1001 established mutational contexts DBS-78 and ID-83<sup>19,57</sup> were used for the extraction of DBS and  
1002 ID signatures (**Supplementary Tables 6-7**).

1003  
1004 Copy number signatures were extracted *de novo* using SigProfilerExtractor with default  
1005 parameters and following an updated context definition benefitting from WGS data (CN-68)  
1006 (**Supplementary Table 8**), which allowed to further characterize CN segments below 100kbp in  
1007 length (in contrast to current COSMICv3.4 reference signatures using the CN-48 context, which  
1008 were based on SNP6 microarray data and therefore without the resolution to characterize short CN  
1009 segments)<sup>61</sup>. SV signatures were extracted using a similarly refined context, with an in-depth  
1010 characterization of short SV alterations below 1kbp (SV-38 context, in contrast to current  
1011 COSMICv3.4 signatures based on the SV-32 context<sup>62</sup>; **Supplementary Table 9**).

1012  
1013 After *de novo* extraction was completed, SigProfilerAssignment<sup>42</sup> v0.0.29 was used to decompose  
1014 the *de novo* extracted SBS, ID, DBS, CN, and SV mutational signatures into COSMICv3.4  
1015 reference signatures based on the GRCh38 reference genome<sup>63</sup> (**Supplementary Table 10**). When  
1016 possible, SigProfilerAssignment matched each *de novo* extracted mutational signature to a set of  
1017 previously identified COSMICv3.4 signatures (**Supplementary Note**). For the SBS-288, CN-68,  
1018 and SV-38 signatures, this required collapsing the high-definition classifications into the standard  
1019 SBS-96, CN-48, and SV-32 mutational classifications, respectively. Four of the *de novo* extracted  
1020 MSS SBS signatures did not match any previous COSMICv3.4 signatures, with three of them  
1021 (SBS\_F, SBS\_H, and SBS\_M) showing a strong similarity with previously reported signatures in  
1022 the UK population<sup>21</sup> (cosine similarity > 0.93), and one (SBS\_O) reflecting a cleaner version of a  
1023 previously reported COSMICv3.4 signature (SBS41). To validate the latter, we performed a

1024 decomposition of the current mutational profile of signature SBS41 using the decomposed  
1025 signatures from our analysis, obtaining a confirmation that SBS41 can be reconstructed by a linear  
1026 combination of SBS\_O (contributing 19.00% of the mutational profile), SBS93 (62.54%), and  
1027 SBS34 (12.60%), and SBS5 (5.86%) with a cosine similarity of 0.91. Notably, SBS93, first  
1028 identified in gastric tumors<sup>35</sup>, was unknown at the time SBS41 was first reported<sup>19</sup>. For the MSS  
1029 cohort, one ID (ID\_J), one CN (CN\_F), and two SV signatures (SV\_B and SV\_D) were  
1030 additionally not decomposed into previously known signatures, and therefore considered as novel  
1031 (**Supplementary Table 10**). The novel SV signature SV\_D, identified in the MSS cohort, was  
1032 also considered for the decomposition of *de novo* SV signatures extracted in the MSI cohort. In  
1033 the MSI cohort, four of the *de novo* extracted SBS signatures (SBS\_I\_MSI, SBS\_M\_MSI, M,  
1034 SBS\_N\_MSI, and SBS\_O\_MSI) as well as one *de novo* DBS signature (DBS\_B\_MSI) did not  
1035 match COSMICv3.4 signatures, with SBS\_M\_MSI showing a strong similarity with a previously  
1036 reported signature in the UK population<sup>21</sup> (cosine similarity=0.89), and the other four signatures  
1037 considered as novel (**Supplementary Note**).

1038  
1039 mSigHdp<sup>59</sup> extraction of SBS-96 and ID-83 signatures was performed on the 802 MSS subset  
1040 using the suggested parameters and using the country of origin to construct the hierarchy.  
1041 SigProfilerAssignment was subsequently used to match mSigHdp *de novo* signatures to previously  
1042 identified COSMIC signatures. A comparison of the signatures extracted from mSigHdp and  
1043 SigProfilerExtractor can be found in the **Supplementary Note**.

1044

1045 **Attribution of mutational signatures to individual samples**

1046 Known COSMIC signatures and *de novo* signatures that were not decomposed into COSMIC  
1047 signatures (**Supplementary Table 10; Supplementary Note**) were attributed for each sample  
1048 using MSA<sup>64</sup> (v2.0) for SBS, ID, and DBS, whereas SigProfilerAssignment<sup>42</sup> was used for CN and  
1049 SV. A conservative approach was used for MSA attributions utilizing the  
1050 (params.no\_CI\_for\_penalties=False) option for the calculation of optimum penalties. Pruned  
1051 attributions were used for the final analysis, where confidence intervals were applied to each  
1052 attributed mutational signature and any signature activity with a lower confidence limit equal to 0  
1053 was removed.

1054

#### 1055 **Attribution of mutational signatures to individual somatic mutations**

1056 SBS and ID mutational signatures were probabilistically attributed to individual somatic mutations  
1057 using the MSA activities per sample, based on Bayes' rule and the specific mutational context for  
1058 the mutation, as previously described<sup>42</sup>. Briefly, to calculate the probability of a specific mutational  
1059 signature being responsible for a mutation in a given mutational context and in a particular sample,  
1060 we multiplied the general probability of the signature causing mutations in a specific mutational  
1061 context (obtained from the mutational signature profile) by the activity of the signature in the  
1062 sample (obtained from the signature activities), and then normalized this value dividing by the  
1063 total number of mutations corresponding to the specific mutational context (obtained from the  
1064 reconstructed mutational profile of the sample). The signature with the maximum likelihood  
1065 estimation was assigned to each individual somatic mutation.

1066

#### 1067 **Driver gene analysis**

1068 Consensus *de novo* driver gene identification was performed by IntOGen<sup>40</sup>, which combines seven  
1069 state-of-the-art computational methods to detect signals of positive selection across the cohort. The  
1070 genes identified as drivers with a combination q-value<0.10 were classified according to their  
1071 mode of action in tumorigenesis (*i.e.*, tumor suppressor genes or oncogenes) based on the  
1072 relationship between the excess of observed nonsynonymous and truncating mutations computed  
1073 by dNdScv<sup>65</sup> and their annotations in the Cancer Gene Census<sup>66</sup>.

1074  
1075 To identify potential driver mutations, we selected SBS or ID mutations that fulfilled any of the  
1076 following criteria: mutations classified as “Oncogenic” or “Likely Oncogenic” by OncoKB<sup>67</sup>;  
1077 mutations classified as drivers in the TCGA MC3 drivers study<sup>68</sup>; truncating mutations in driver  
1078 genes annotated as tumor suppressors; recurrent missense mutations (seen in at least three cases);  
1079 mutations classified as “Likely Drivers” by boostDM (score >0.50)<sup>69</sup>; or missense mutations  
1080 classified as “Likely Pathogenic” by AlphaMissense<sup>70</sup> in driver genes annotated as tumor  
1081 suppressors. Six of the IntOGen-identified driver genes did not carry any potential driver mutations  
1082 according to our strict criteria and were therefore excluded from subsequent analysis. In summary,  
1083 60 driver genes were identified (46 and 31 for MSS and MSI cases, respectively; **Supplementary**  
1084 **Table 21; Supplementary Note**).

1085  
1086

### 1087 **Evolutionary analysis**

1088 DPCLust<sup>53</sup> was run on all complete pipeline MSS samples with Battenberg data ( $n=774$ ) to identify  
1089 clonal structure in each sample. The DPCLust output was used in running MutationTimeR<sup>38</sup> to  
1090 annotate somatic mutations as early clonal, late clonal, subclonal, or NA clonal. Samples with at



1091 least 256 early clonal and late clonal SBSs or 100 early clonal and late clonal IDs were retained  
1092 and split into separate VCF files ( $n=574$  for SBS;  $n=430$  for ID). MSA<sup>64</sup> was run on the resulting  
1093 VCF files to identify the active mutational signatures in the early clonal and late clonal mutations.  
1094 SBS signatures that were found to generate early clonal SBSs in fewer than 50 samples and also  
1095 generated late clonal SBSs in fewer than 50 samples were excluded from the analysis. Similarly,  
1096 ID signatures generating early clonal IDs in fewer than 20 samples and late clonal IDs in fewer  
1097 than 20 samples were also excluded. Wilcoxon signed-rank tests were used to assess the  
1098 differences in the relative activity of each signature between the early clonal and late clonal  
1099 mutations. P-values were adjusted across signatures using the Benjamini-Hochberg method<sup>71</sup>, and  
1100 adjusted p-values were reported as q-values. This process was repeated with the same thresholds  
1101 for SBSs and IDs to also assess the difference in the relative activity of each signature between  
1102 clonal and subclonal mutations ( $n=133$  for SBS;  $n=64$  for ID). Due to the lower numbers,  
1103 signatures that were found to generate clonal somatic mutations in fewer than 10 samples and also  
1104 generated subclonal somatic mutations in fewer than 10 samples were excluded from the analysis.

1105

#### 1106 **Motif analysis**

1107 MutaGene<sup>72</sup> was used to find the number of mutations with the WAWW[T>N]W motif, previously  
1108 associated with colibactin mutagenesis<sup>33</sup>, in each sample, regardless of the DNA strand. This value  
1109 was then divided by the total number of W[T>N]W mutations per sample to identify the percentage  
1110 of W[T>N]W mutations with the colibactin mutational motif.

1111

#### 1112 **Microbiome analysis**

1113 To identify microbial reads that map to the *pks* island (*pks*), non-human reads were aligned to the  
1114 IHE3034 genome (RefSeq assembly: GCF\_000025745.1) using Bowtie2<sup>73</sup>. IHE3034 is a *pks E. coli*  
1115 strain that contains the *pks* island with all 19 *clb* genes in the *clbA-clbS* gene cluster. Prior to  
1116 alignment, poor quality reads were filtered using fastp<sup>74</sup>, and the remaining human reads were  
1117 removed by excluding those that mapped to GRCh38, T2T-CHM13v2.0, and the 47 pangenomes<sup>75</sup>.  
1118 A sample was considered *pks+* if it had at least one read across at least 8 out of the 19 genes in the  
1119 *clbA-clbS* gene cluster. Genome coverage circos plots were generated using Reads Per Kilobase  
1120 per Million (RPKM) values and visualized with the *circlize* R package<sup>76</sup>.

1121

## 1122 **Regressions**

1123 To compare the mutation burden of different variant types, a linear regression of the mutation  
1124 burden logarithm (base 10) was considered, using age, sex, tumor subsite, country, and tumor  
1125 purity as independent variables. For mutational signature-based analyses, signature attributions  
1126 were dichotomized into presence and absence using confidence intervals, with presence defined  
1127 as both lower and upper limits being positive and absence as the lower limit being zero. If a  
1128 signature was present in at least 70% of cases (SBS1, SBS5, SBS18, ID1, ID2, ID14 and CN2 for  
1129 MSS cases; ID1, ID2, DBS\_B\_MSI, CN1, and SV\_D for MSI cases), it was dichotomized into  
1130 above and below the median of attributed mutation counts. The binary attributions served as  
1131 dependent variables in logistic regressions. Regressions with variables presenting complete or  
1132 quasi-complete separation<sup>77</sup> were performed using Firth's bias-reduced logistic regressions based  
1133 on the *logistf* R package. To adjust for confounding factors, sex, age of diagnosis, tumor subsite,  
1134 country, and tumor purity were added as covariates in all regressions, serving as independent  
1135 variables for the regressions. The tumor subsite variable was categorized as proximal colon (ICD-

1136 10-CM codes C18.0, C18.2, C18.3, and C18.4), distal colon (C18.5, C18.6, and C18.7), or rectum  
1137 (C19 and C20), unless otherwise specified. One MSI tumor from an unspecified subsite was  
1138 removed for the multivariable regression models in MSI cases. The age of diagnosis variable was  
1139 generally considered as a numerical variable, or categorized into two (early-onset, <50 years old;  
1140 and late-onset,  $\geq 50$ ) or five subgroups (0-39, 40-49, 50-59, 60-69,  $\geq 70$ ), depending on the analysis  
1141 performed, with specific indications in the corresponding figure legends. Similarly, regressions for  
1142 driver mutations in cancer driver genes and hotspot driver mutations (present in at least 10 cases)  
1143 were done using the same logistic regression models but replacing signature by driver mutation  
1144 prevalence across samples.

1145  
1146 Regressions with colorectal cancer incidence were performed as linear regressions with signature  
1147 attributions with confidence intervals not consistent with zero as dependent variables, and age-  
1148 standardized rates (ASR) of colorectal cancer (and independent ASR of colon and rectal cancer)  
1149 obtained from the Global Cancer Observatory (GLOBOCAN)<sup>1</sup>, sex, age of diagnosis, tumor  
1150 subsite, and tumor purity as independent variables. Regressions were performed on a sample basis.

1151  
1152 Regressions with colibactin presence (based on genomic and/or microbiome-derived detection)  
1153 were performed as linear regressions with age of diagnosis as the dependent variable, and sex,  
1154 tumor subsite, country, and tumor purity as independent variables.

### 1155 **Additional statistical analyses**

1156 For regressions of signatures, driver mutations in cancer driver genes, and hotspot driver  
1157 mutations, p-values were adjusted for multiple comparisons based on the total number of  
1158 decomposed reference mutational signatures considered per variant type (*i.e.*, 19 SBS, 7 DBS, 11

1159 ID, 9 CN, and 11 SV signatures for MSS cases; 18 SBS, 10 DBS, 2 ID, 4 CN, and 4 SV for MSI  
1160 cases), cancer genes (46 for MSS; 31 for MSI), or hotspot driver mutations (38 for MSS; 14 for  
1161 MSI) using the Benjamini-Hochberg method<sup>71</sup>. For country enrichment analyses, the mutation  
1162 burdens and binary attributions of mutational signatures were compared for each country against  
1163 all others. Therefore, p-values were also adjusted for multiple comparisons based on the total  
1164 number of countries assessed (a total of 11 countries). Adjusted p-values were reported as q-values,  
1165 with q-values<0.05 considered statistically significant. For age of diagnosis-based regressions of  
1166 colibactin presence across tumor subsites, p-values were adjusted and reported as q-values based  
1167 on the total number of tumor subsites assessed (a total of 3 tumor subsites). For the age of diagnosis  
1168 trend enrichment analysis of signatures, p-trends were reported, with p-trends<0.05 considered  
1169 statistically significant. For evidence of co-occurrence or mutual exclusivity of two signatures,  
1170 two-sided Fisher's exact tests were used, and p-values were reported, with p-values<0.05  
1171 considered statistically significant.

1172

1173 **DATA AVAILABILITY**

1174 Whole-genome sequencing data, somatic mutations, and patient metadata are deposited in the  
1175 European Genome-phenome Archive (EGA) associated with study EGAS00001003774. All other  
1176 data is provided in the accompanying Supplementary Tables.

1177

1178 **CODE AVAILABILITY**

1179 All algorithms used for data analysis are publicly available with repositories noted within the  
1180 respective method sections. The code used for regression analysis and figures is available at  
1181 [https://github.com/AlexandrovLab/Mutographs\\_CRC](https://github.com/AlexandrovLab/Mutographs_CRC).

1182

1183 **ACKNOWLEDGMENTS**

1184 The authors thank the IARC General Services, including the Laboratory Services and Biobank  
1185 team led by Z. Kozlakidis and the Section of Support to Research overseen by C. Mehta under  
1186 IARC regular budget funding for the support provided; Laura O'Neill, Kirsty Roberts, Katie  
1187 Smith, Siobhan Austin-Guest, and the staff of Sequencing Operations at the Wellcome Sanger  
1188 Institute for their contribution; Laura Rodríguez Porras for her help in designing and reviewing the  
1189 figures; the work of all other collaborators in the Mutographs project who participated in the  
1190 recruitment of patients in all centers; and all the patients involved in this study and their families.  
1191 The computational analyses reported in this manuscript have utilized the Triton Shared Computing  
1192 Cluster at the San Diego Supercomputer Center of UC San Diego. Where authors are identified as  
1193 personnel of the International Agency for Research on Cancer / World Health Organization, the  
1194 authors alone are responsible for the views expressed in this article and they do not necessarily

1195 represent the decisions, policy or views of the International Agency for Research on Cancer /  
1196 World Health Organization.

1197

## 1198 **FUNDING**

1199 This work was delivered as part of the Mutographs team supported by the Cancer Grand  
1200 Challenges partnership funded by Cancer Research UK (C98/A24032). Work at UC San Diego  
1201 was also supported by the US National Institute of Health (NIH) grants R01ES032547-01,  
1202 R01CA269919-01, and 1U01CA290479-01 to L.B.A. as well as by L.B.A.'s Packard Fellowship  
1203 for Science and Engineering. The research performed in L.B.A.'s lab was further supported by UC  
1204 San Diego Sanford Stem Cell Institute. This work was supported in part by an IARC Fellowship  
1205 Award to Wellington Oliveira dos Santos through The Mark Foundation for Cancer Research.  
1206 Work at the IARC/WHO was also supported by regular budget funding. Work at the Wellcome  
1207 Sanger Institute was also supported by the Wellcome Trust (grants 206194 and 220540/Z/20/A).  
1208 Work at Masaryk Memorial Cancer Institute, Brno, Czech Republic, was supported by MH CZ -  
1209 DRO (MMCI, 00209805). Porto Alegre center in Brazil received support from Hospital de Clínicas  
1210 de Porto Alegre and Fundação Médica do Rio Grande do Sul. Barretos Cancer Hospital, in Brazil,  
1211 was also supported by the Public Ministry of Labor Campinas (Research, Prevention, and  
1212 Education of Occupational Cancer). This work was supported by grants from Practical Research  
1213 for Innovative Cancer Control from the Japan Agency for Medical Research and Development  
1214 (AMED) (JP 24ck0106800h0002 to T.S.) and the National Cancer Center Research and  
1215 Development Fund (2023-A-05 to T.S.). Work at Sinai Health System, Toronto, Canada received  
1216 support from the NIH (grant U01CA167551). The funders had no roles in study design, data  
1217 collection and analysis, decision to publish, or preparation of the manuscript.

1218 **AUTHOR CONTRIBUTIONS**

1219 The study was conceived, designed and supervised by M.R.S., P.B., and L.B.A. Analysis of data  
1220 was performed by M.D.-G., W.d.S., S. Moody, M.K., A.A., C.D.S, R.V., S. Senkin, J.W., S.F.,  
1221 E.N.B., A.K., B.O., T. Cattiaux, R.C.C.P., V.G., S.C., and J.W.T. Analysis and interpretation of  
1222 the microbiomics data was performed by A.A. with assistance and advice from L.B.A. and T.D.L.  
1223 Pathology review was carried out by B.A.-A. Sample manipulation was carried out by P.C., C.C.,  
1224 and C.L. Patient and sample recruitment was led or facilitated by D.Z., R.C., M.A., L.P., S.G.,  
1225 J.Y., R.M., A.N., M.M., K.E., S. Milosavljevic, S. Sangrajrang, M.P.C., S.A., R.M.R., M.T.R.,  
1226 L.G.R., D.P.G., I.H., J.K., C.A.V., T.A.P., B.Ś., J.L., K.R.P., A.H.-S., T.S., S. Shiba, S.  
1227 Sangkhathat, T. Chitapanarux, G.R., P.A.-P., D.C.D., and F.H.d.O. Scientific project management  
1228 was carried out by L.H., A.C.d.C., and S.P. M.D.-G., W.d.S., and S. Moody jointly contributed  
1229 and were responsible for overall scientific coordination. The manuscript was written by M.D.-G.,  
1230 W.d.S., S. Moody, M.R.S., P.B., and L.B.A., with contributions from all other authors. All authors  
1231 read and approved the final manuscript.

1232

1233 **COMPETING INTERESTS**

1234 L.B.A. is a co-founder, CSO, scientific advisory member, and consultant for io9, has equity and  
1235 receives income. The terms of this arrangement have been reviewed and approved by the  
1236 University of California, San Diego in accordance with its conflict of interest policies. L.B.A. is  
1237 also a compensated member of the scientific advisory board of Inocras. L.B.A.'s spouse is an  
1238 employee of Hologic, Inc. E.N.B. is a consultant for io9, has equity, and receives income. A.A.  
1239 and L.B.A. declare U.S. provisional patent application filed with UCSD with serial number  
1240 63/366,392. E.N.B. and L.B.A. declare U.S. provisional patent application filed with UCSD with

1241 serial numbers 63/269,033. L.B.A. also declares U.S. provisional applications filed with UCSD  
1242 with serial numbers: 63/289,601; 63/412,835; as well as an international patent application  
1243 PCT/US2023/010679. L.B.A. is also an inventor of a US Patent 10,776,718 for source  
1244 identification by non-negative matrix factorization. M.R.S. is founder, consultant, and stockholder  
1245 for Quotient Therapeutics. L.B.A., M.D.-G., P.B., S.P., M.R.S., and S. Moody declare a European  
1246 patent application with application number EP25305077.7. T.D.L. is a co-founder and CSO of  
1247 Microbiotica. All other authors declare that they have no competing interests.  
1248



## 1249 REFERENCES

- 1250 1 Bray, F. *et al.* Global cancer statistics 2022: GLOBOCAN estimates of incidence and  
1251 mortality worldwide for 36 cancers in 185 countries. *CA Cancer J Clin* **74**, 229-263  
1252 (2024). <https://doi.org/10.3322/caac.21834>
- 1253 2 Brennan, P. & Davey-Smith, G. Identifying Novel Causes of Cancers to Enhance Cancer  
1254 Prevention: New Strategies Are Needed. *J Natl Cancer Inst* **114**, 353-360 (2022).  
1255 <https://doi.org/10.1093/jnci/djab204>
- 1256 3 Kucab, J. E. *et al.* A Compendium of Mutational Signatures of Environmental Agents.  
1257 *Cell* **177**, 821-836 e816 (2019). <https://doi.org/10.1016/j.cell.2019.03.001>
- 1258 4 Ames, B. N., Durston, W. E., Yamasaki, E. & Lee, F. D. Carcinogens are mutagens: a  
1259 simple test system combining liver homogenates for activation and bacteria for detection.  
1260 *Proc Natl Acad Sci U S A* **70**, 2281-2285 (1973). <https://doi.org/10.1073/pnas.70.8.2281>
- 1261 5 Senkin, S. *et al.* Geographic variation of mutagenic exposures in kidney cancer genomes.  
1262 *Nature* (2024). <https://doi.org/10.1038/s41586-024-07368-2>
- 1263 6 Moody, S. *et al.* Mutational signatures in esophageal squamous cell carcinoma from eight  
1264 countries with varying incidence. *Nat Genet* **53**, 1553-1563 (2021).  
1265 <https://doi.org/10.1038/s41588-021-00928-6>
- 1266 7 Zhang, T. *et al.* Genomic and evolutionary classification of lung cancer in never smokers.  
1267 *Nature Genetics* **53**, 1348-1359 (2021). <https://doi.org/10.1038/s41588-021-00920-0>
- 1268 8 Perdomo, S. *et al.* The Mutographs biorepository: A unique genomic resource to study  
1269 cancer around the world. *Cell Genom* **4**, 100500 (2024).  
1270 <https://doi.org/10.1016/j.xgen.2024.100500>
- 1271 9 Torrens, L. *et al.* The Complexity of Tobacco Smoke-Induced Mutagenesis in Head and  
1272 Neck Cancer. *medRxiv*, 2024.2004.2015.24305006 (2024).  
1273 <https://doi.org/10.1101/2024.04.15.24305006>
- 1274 10 Patel, S. G., Karlitz, J. J., Yen, T., Lieu, C. H. & Boland, C. R. The rising tide of early-  
1275 onset colorectal cancer: a comprehensive review of epidemiology, clinical features,  
1276 biology, risk factors, prevention, and early detection. *Lancet Gastroenterol Hepatol* **7**,  
1277 262-274 (2022). [https://doi.org/10.1016/S2468-1253\(21\)00426-X](https://doi.org/10.1016/S2468-1253(21)00426-X)
- 1278 11 Siegel, R. L. *et al.* Global patterns and trends in colorectal cancer incidence in young  
1279 adults. *Gut* **68**, 2179-2185 (2019). <https://doi.org/10.1136/gutjnl-2019-319511>
- 1280 12 Siegel, R. L., Jemal, A. & Ward, E. M. Increase in incidence of colorectal cancer among  
1281 young men and women in the United States. *Cancer Epidemiol Biomarkers Prev* **18**,  
1282 1695-1698 (2009). <https://doi.org/10.1158/1055-9965.EPI-09-0186>
- 1283 13 Sinicrope, F. A. Increasing Incidence of Early-Onset Colorectal Cancer. *N Engl J Med*  
1284 **386**, 1547-1558 (2022). <https://doi.org/10.1056/NEJMra2200869>
- 1285 14 Vuik, F. E. *et al.* Increasing incidence of colorectal cancer in young adults in Europe over  
1286 the last 25 years. *Gut* **68**, 1820-1826 (2019). <https://doi.org/10.1136/gutjnl-2018-317592>
- 1287 15 Spaander, M. C. W. *et al.* Young-onset colorectal cancer. *Nat Rev Dis Primers* **9**, 21  
1288 (2023). <https://doi.org/10.1038/s41572-023-00432-7>
- 1289 16 Stigliano, V., Sanchez-Mete, L., Martayan, A. & Anti, M. Early-onset colorectal cancer: a  
1290 sporadic or inherited disease? *World J Gastroenterol* **20**, 12420-12430 (2014).  
1291 <https://doi.org/10.3748/wjg.v20.i35.12420>

- 1292 17 You, Y. N., Xing, Y., Feig, B. W., Chang, G. J. & Cormier, J. N. Young-onset colorectal  
1293 cancer: is it time to pay attention? *Arch Intern Med* **172**, 287-289 (2012).  
1294 <https://doi.org/10.1001/archinternmed.2011.602>
- 1295 18 Venugopal, A. & Carethers, J. M. Epidemiology and biology of early onset colorectal  
1296 cancer. *EXCLI J* **21**, 162-182 (2022). <https://doi.org/10.17179/excli2021-4456>
- 1297 19 Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature*  
1298 **578**, 94-101 (2020). <https://doi.org/10.1038/s41586-020-1943-3>
- 1299 20 Cancer Genome Atlas, N. *et al.* Comprehensive molecular characterization of human  
1300 colon and rectal cancer. *Nature* **487**, 330-337 (2012). <https://doi.org/10.1038/nature11252>
- 1301 21 Degasperi, A. *et al.* Substitution mutational signatures in whole-genome-sequenced  
1302 cancers in the UK population. *Science* **376**, abl9283 (2022).  
1303 <https://doi.org/10.1126/science.abl9283>
- 1304 22 Cornish, A. J. *et al.* The genomic landscape of 2,023 colorectal cancers. *Nature* (2024).  
1305 <https://doi.org/10.1038/s41586-024-07747-9>
- 1306 23 Priestley, P. *et al.* Pan-cancer whole-genome analyses of metastatic solid tumours. *Nature*  
1307 **575**, 210-216 (2019). <https://doi.org/10.1038/s41586-019-1689-y>
- 1308 24 Martinez-Jimenez, F. *et al.* Pan-cancer whole-genome comparison of primary and  
1309 metastatic solid tumours. *Nature* **618**, 333-341 (2023). <https://doi.org/10.1038/s41586-023-06054-z>
- 1310
- 1311 25 Mendelaar, P. A. J. *et al.* Whole genome sequencing of metastatic colorectal cancer  
1312 reveals prior treatment effects and specific metastasis features. *Nat Commun* **12**, 574  
1313 (2021). <https://doi.org/10.1038/s41467-020-20887-6>
- 1314 26 Rosendahl Huber, A. *et al.* Improved detection of colibactin-induced mutations by  
1315 genotoxic E. coli in organoids and colorectal cancer. *Cancer Cell* **42**, 487-496 (2024).  
1316 <https://doi.org/10.1016/j.ccell.2024.02.009>
- 1317 27 Nunes, L. *et al.* Prognostic genome and transcriptome signatures in colorectal cancers.  
1318 *Nature* (2024). <https://doi.org/10.1038/s41586-024-07769-3>
- 1319 28 Díaz-Gay, M. & Alexandrov, L. B. in *Advances in Cancer Research* Vol. 151 (eds  
1320 Franklin G. Berger & C. Richard Boland) 385-424 (Academic Press, 2021).
- 1321 29 Helleday, T., Eshtad, S. & Nik-Zainal, S. Mechanisms underlying mutational signatures  
1322 in human cancers. *Nat Rev Genet* **15**, 585-598 (2014). <https://doi.org/10.1038/nrg3729>
- 1323 30 Pich, O. *et al.* The mutational footprints of cancer therapies. *Nat Genet* **51**, 1732-1740  
1324 (2019). <https://doi.org/10.1038/s41588-019-0525-5>
- 1325 31 Nik-Zainal, S. *et al.* Mutational processes molding the genomes of 21 breast cancers. *Cell*  
1326 **149**, 979-993 (2012). <https://doi.org/10.1016/j.cell.2012.04.024>
- 1327 32 Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**,  
1328 415-421 (2013). <https://doi.org/10.1038/nature12477>
- 1329 33 Pleguezuelos-Manzano, C. *et al.* Mutational signature in colorectal cancer caused by  
1330 genotoxic pks(+) E. coli. *Nature* **580**, 269-273 (2020). <https://doi.org/10.1038/s41586-020-2080-8>
- 1331
- 1332 34 Lee-Six, H. *et al.* The landscape of somatic mutation in normal colorectal epithelial cells.  
1333 *Nature* **574**, 532-537 (2019). <https://doi.org/10.1038/s41586-019-1672-7>
- 1334 35 Islam, S. M. A. *et al.* Uncovering novel mutational signatures by de novo extraction with  
1335 SigProfilerExtractor. *Cell Genom* **2**, 100179 (2022).  
1336 <https://doi.org/10.1016/j.xgen.2022.100179>

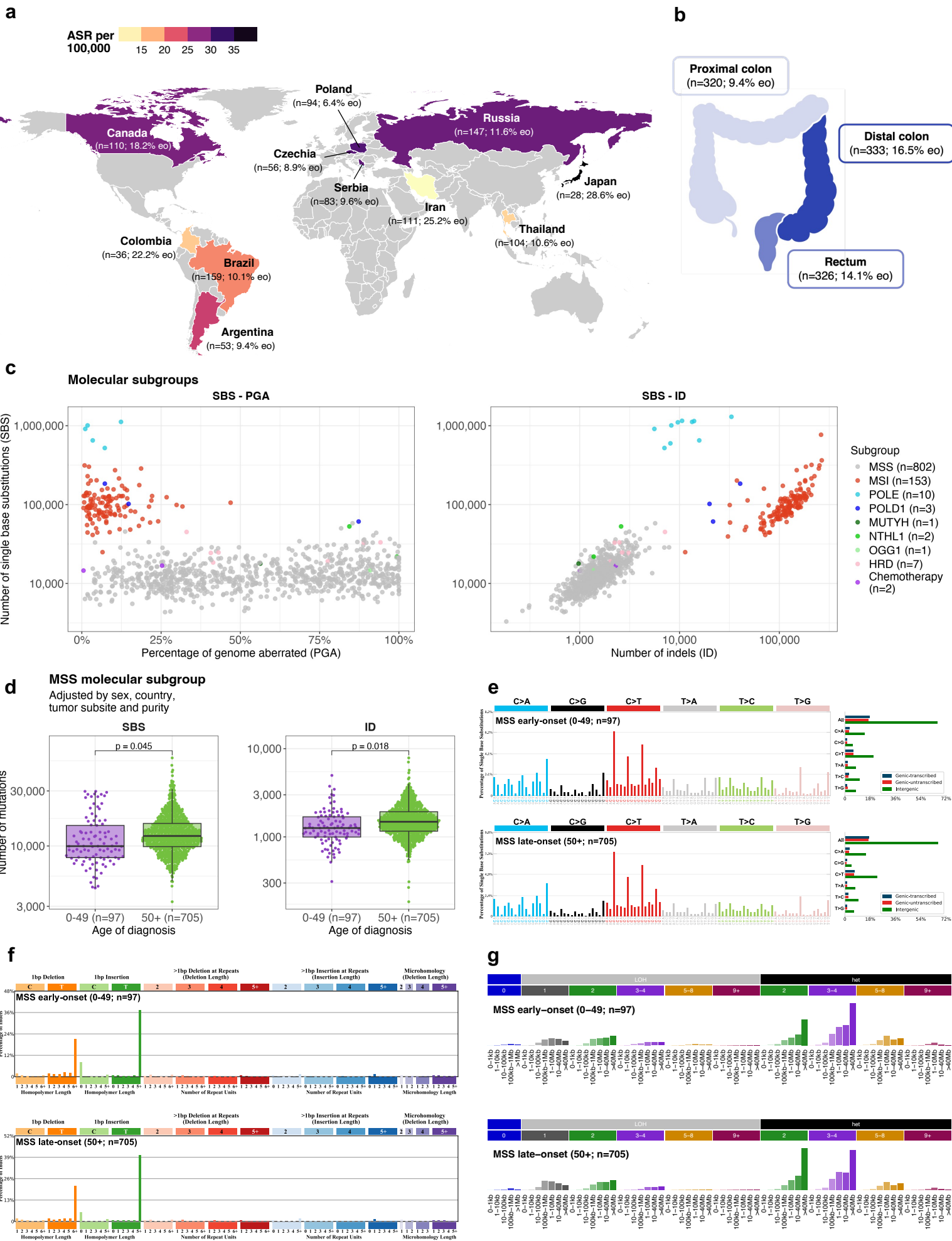
- 1337 36 Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells. *Nat*  
1338 *Genet* **47**, 1402-1407 (2015). <https://doi.org/10.1038/ng.3441>
- 1339 37 D'Entropio, S. C. *et al.* Characterizing genetic intra-tumor heterogeneity across 2,658 human  
1340 cancer genomes. *Cell* **184**, 2239-2254 (2021). <https://doi.org/10.1016/j.cell.2021.03.009>
- 1341 38 Gerstung, M. *et al.* The evolutionary history of 2,658 cancers. *Nature* **578**, 122-128  
1342 (2020). <https://doi.org/10.1038/s41586-019-1907-7>
- 1343 39 Chen, B. *et al.* Contribution of pks(+) *E. coli* mutations to colorectal carcinogenesis. *Nat*  
1344 *Commun* **14**, 7827 (2023). <https://doi.org/10.1038/s41467-023-43329-5>
- 1345 40 Martinez-Jimenez, F. *et al.* A compendium of mutational cancer driver genes. *Nat Rev*  
1346 *Cancer* **20**, 555-572 (2020). <https://doi.org/10.1038/s41568-020-0290-x>
- 1347 41 Kim, J. E. *et al.* High prevalence of TP53 loss and whole-genome doubling in early-onset  
1348 colorectal cancer. *Exp Mol Med* **53**, 446-456 (2021). [https://doi.org/10.1038/s12276-021-](https://doi.org/10.1038/s12276-021-00583-1)  
1349 [00583-1](https://doi.org/10.1038/s12276-021-00583-1)
- 1350 42 Díaz-Gay, M. *et al.* Assigning mutational signatures to individual samples and individual  
1351 somatic mutations with SigProfilerAssignment. *Bioinformatics* **39**, btad756 (2023).  
1352 <https://doi.org/10.1093/bioinformatics/btad756>
- 1353 43 Terlouw, D. *et al.* Recurrent APC Splice Variant c.835-8A>G in Patients With  
1354 Unexplained Colorectal Polyposis Fulfilling the Colibactin Mutational Signature.  
1355 *Gastroenterology* **159**, 1612-1614 e1615 (2020).  
1356 <https://doi.org/10.1053/j.gastro.2020.06.055>
- 1357 44 Fearon, E. R. & Vogelstein, B. A genetic model for colorectal tumorigenesis. *Cell* **61**,  
1358 759-767 (1990). [https://doi.org/10.1016/0092-8674\(90\)90186-i](https://doi.org/10.1016/0092-8674(90)90186-i)
- 1359 45 Carethers, J. M. & Jung, B. H. Genetics and Genetic Biomarkers in Sporadic Colorectal  
1360 Cancer. *Gastroenterology* **149**, 1177-1190 e1173 (2015).  
1361 <https://doi.org/10.1053/j.gastro.2015.06.047>
- 1362 46 Perdomo, S. Mutational signatures in five cancer types across five continents. Standard  
1363 Operating Procedures (SOPs). *Zenodo* (2024). <https://doi.org/10.5281/zenodo.11836372>
- 1364 47 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler  
1365 transform. *Bioinformatics* **25**, 1754-1760 (2009).  
1366 <https://doi.org/10.1093/bioinformatics/btp324>
- 1367 48 Whalley, J. P. *et al.* Framework for quality assessment of whole genome cancer  
1368 sequences. *Nat Commun* **11**, 5040 (2020). <https://doi.org/10.1038/s41467-020-18688-y>
- 1369 49 Bergmann, E. A., Chen, B. J., Arora, K., Vacic, V. & Zody, M. C. Conpair: concordance  
1370 and contamination estimator for matched tumor-normal pairs. *Bioinformatics* **32**, 3196-  
1371 3198 (2016). <https://doi.org/10.1093/bioinformatics/btw389>
- 1372 50 Kim, S. *et al.* Strelka2: fast and accurate calling of germline and somatic variants. *Nat*  
1373 *Methods* **15**, 591-594 (2018). <https://doi.org/10.1038/s41592-018-0051-x>
- 1374 51 Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic  
1375 variants from high-throughput sequencing data. *Nucleic Acids Res* **38**, e164 (2010).  
1376 <https://doi.org/10.1093/nar/gkq603>
- 1377 52 Van Loo, P. *et al.* Allele-specific copy number analysis of tumors. *Proc Natl Acad Sci U S*  
1378 *A* **107**, 16910-16915 (2010). <https://doi.org/10.1073/pnas.1009843107>
- 1379 53 Nik-Zainal, S. *et al.* The life history of 21 breast cancers. *Cell* **149**, 994-1007 (2012).  
1380 <https://doi.org/10.1016/j.cell.2012.04.023>

- 1381 54 Jones, D. *et al.* cgpCaVEManWrapper: Simple Execution of CaVEMan in Order to  
1382 Detect Somatic Single Nucleotide Variants in NGS Data. *Curr Protoc Bioinformatics* **56**,  
1383 15 10 11-15 10 18 (2016). <https://doi.org/10.1002/cpbi.20>
- 1384 55 Raine, K. M. *et al.* cgpPindel: Identifying Somatically Acquired Insertion and Deletion  
1385 Events from Paired End Sequencing. *Curr Protoc Bioinformatics* **52**, 15 17 11-15 17 12  
1386 (2015). <https://doi.org/10.1002/0471250953.bi1507s52>
- 1387 56 Khandekar, A. *et al.* Visualizing and exploring patterns of large mutational events with  
1388 SigProfilerMatrixGenerator. *BMC Genomics* **24**, 469 (2023).  
1389 <https://doi.org/10.1186/s12864-023-09584-y>
- 1390 57 Bergstrom, E. N. *et al.* SigProfilerMatrixGenerator: a tool for visualizing and exploring  
1391 patterns of small mutational events. *BMC Genomics* **20**, 685 (2019).  
1392 <https://doi.org/10.1186/s12864-019-6041-2>
- 1393 58 Gilson, P. *et al.* Evaluation of 3 molecular-based assays for microsatellite instability  
1394 detection in formalin-fixed tissues of patients with endometrial and colorectal cancers.  
1395 *Sci Rep* **10**, 16386 (2020). <https://doi.org/10.1038/s41598-020-73421-5>
- 1396 59 Liu, M., Wu, Y., Jiang, N., Boot, A. & Rozen, S. G. mSigHdp: hierarchical Dirichlet  
1397 process mixture modeling for mutational signature discovery. *NAR Genom Bioinform* **5**,  
1398 lqad005 (2023). <https://doi.org/10.1093/nargab/lqad005>
- 1399 60 Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R.  
1400 Deciphering signatures of mutational processes operative in human cancer. *Cell Rep* **3**,  
1401 246-259 (2013). <https://doi.org/10.1016/j.celrep.2012.12.008>
- 1402 61 Steele, C. D. *et al.* Signatures of copy number alterations in human cancer. *Nature* **606**,  
1403 984-991 (2022). <https://doi.org/10.1038/s41586-022-04738-6>
- 1404 62 Everall, A. *et al.* Comprehensive repertoire of the chromosomal alteration and mutational  
1405 signatures across 16 cancer types from 10,983 cancer patients. *medRxiv*,  
1406 2023.2006.2007.23290970 (2023). <https://doi.org/10.1101/2023.06.07.23290970>
- 1407 63 Sondka, Z. *et al.* COSMIC: a curated database of somatic variants and clinical data for  
1408 cancer. *Nucleic Acids Res* **52**, D1210-D1217 (2024). <https://doi.org/10.1093/nar/gkad986>
- 1409 64 Senkin, S. MSA: reproducible mutational signature attribution with confidence based on  
1410 simulations. *BMC Bioinformatics* **22**, 540 (2021). <https://doi.org/10.1186/s12859-021-04450-8>
- 1411
- 1412 65 Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues.  
1413 *Cell* **171**, 1029-1041 (2017). <https://doi.org/10.1016/j.cell.2017.09.042>
- 1414 66 Sondka, Z. *et al.* The COSMIC Cancer Gene Census: describing genetic dysfunction  
1415 across all human cancers. *Nat Rev Cancer* **18**, 696-705 (2018).  
1416 <https://doi.org/10.1038/s41568-018-0060-1>
- 1417 67 Chakravarty, D. *et al.* OncoKB: A Precision Oncology Knowledge Base. *JCO Precis*  
1418 *Oncol* **2017**, 1-16 (2017). <https://doi.org/10.1200/PO.17.00011>
- 1419 68 Bailey, M. H. *et al.* Comprehensive Characterization of Cancer Driver Genes and  
1420 Mutations. *Cell* **173**, 371-385 e318 (2018). <https://doi.org/10.1016/j.cell.2018.02.060>
- 1421 69 Muiños, F., Martinez-Jimenez, F., Pich, O., Gonzalez-Perez, A. & Lopez-Bigas, N. In  
1422 silico saturation mutagenesis of cancer genes. *Nature* **596**, 428-432 (2021).  
1423 <https://doi.org/10.1038/s41586-021-03771-1>
- 1424 70 Cheng, J. *et al.* Accurate proteome-wide missense variant effect prediction with  
1425 AlphaMissense. *Science* **381**, eadg7492 (2023). <https://doi.org/10.1126/science.adg7492>

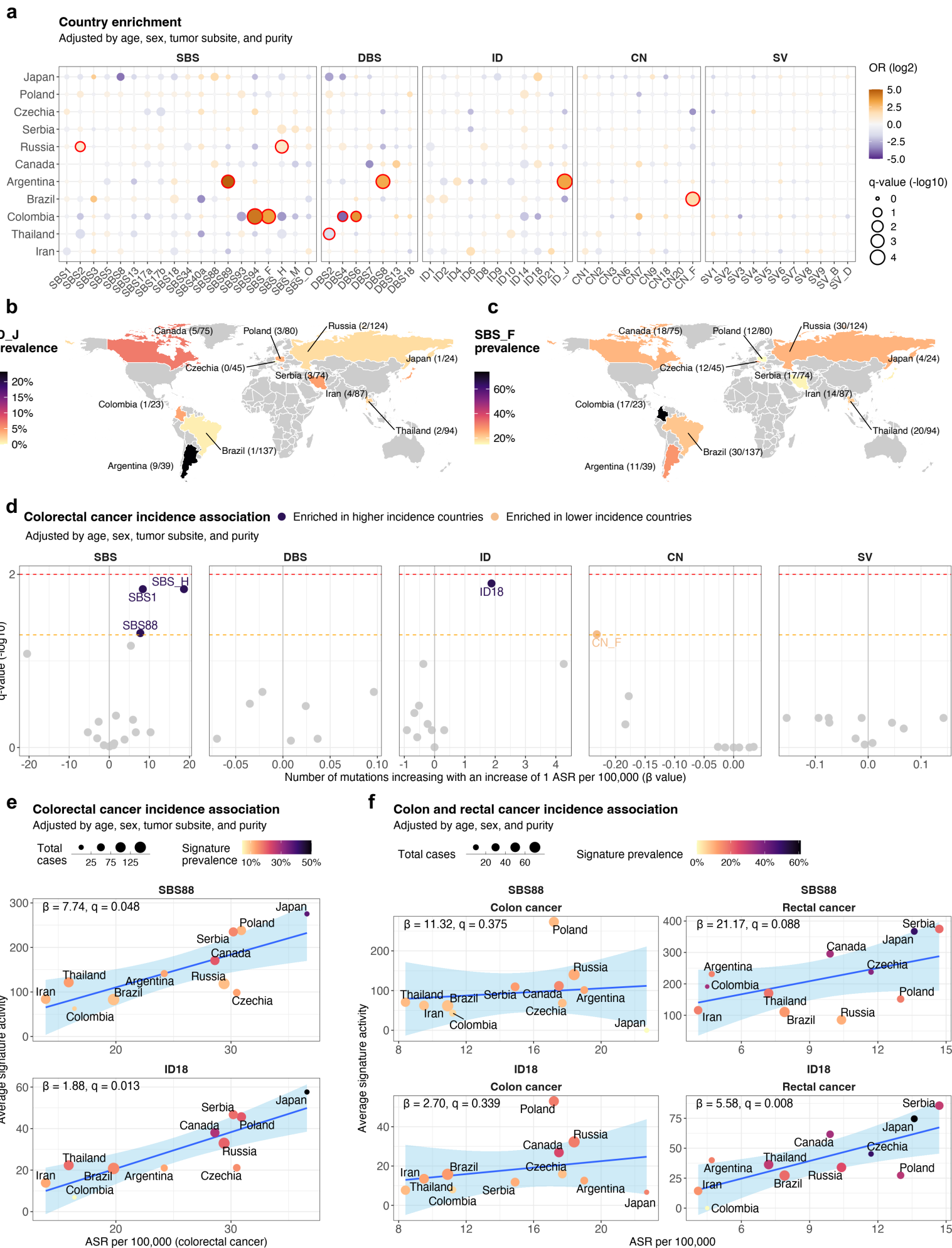


- 1426 71 Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate - a Practical and  
1427 Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B-*  
1428 *Statistical Methodology* **57**, 289-300 (1995). [https://doi.org/10.1111/j.2517-](https://doi.org/10.1111/j.2517-6161.1995.tb02031.x)  
1429 [6161.1995.tb02031.x](https://doi.org/10.1111/j.2517-6161.1995.tb02031.x)
- 1430 72 Goncarenco, A. *et al.* Exploring background mutational processes to decipher cancer  
1431 genetic heterogeneity. *Nucleic Acids Res* **45**, W514-W522 (2017).  
1432 <https://doi.org/10.1093/nar/gkx367>
- 1433 73 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods*  
1434 **9**, 357-359 (2012). <https://doi.org/10.1038/nmeth.1923>
- 1435 74 Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.  
1436 *Bioinformatics* **34**, i884-i890 (2018). <https://doi.org/10.1093/bioinformatics/bty560>
- 1437 75 Liao, W. W. *et al.* A draft human pangenome reference. *Nature* **617**, 312-324 (2023).  
1438 <https://doi.org/10.1038/s41586-023-05896-x>
- 1439 76 Gu, Z., Gu, L., Eils, R., Schlesner, M. & Brors, B. circlize Implements and enhances  
1440 circular visualization in R. *Bioinformatics* **30**, 2811-2812 (2014).  
1441 <https://doi.org/10.1093/bioinformatics/btu393>
- 1442 77 Mansournia, M. A., Geroldinger, A., Greenland, S. & Heinze, G. Separation in Logistic  
1443 Regression: Causes, Consequences, and Control. *Am J Epidemiol* **187**, 864-870 (2018).  
1444 <https://doi.org/10.1093/aje/kwx299>  
1445

# Figure 1



## Figure 2

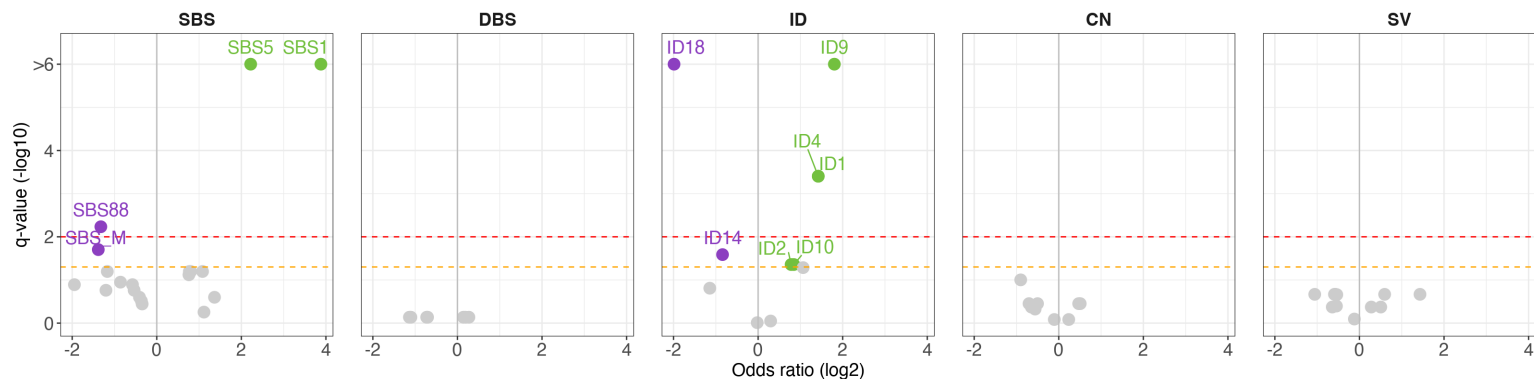


### Figure 3

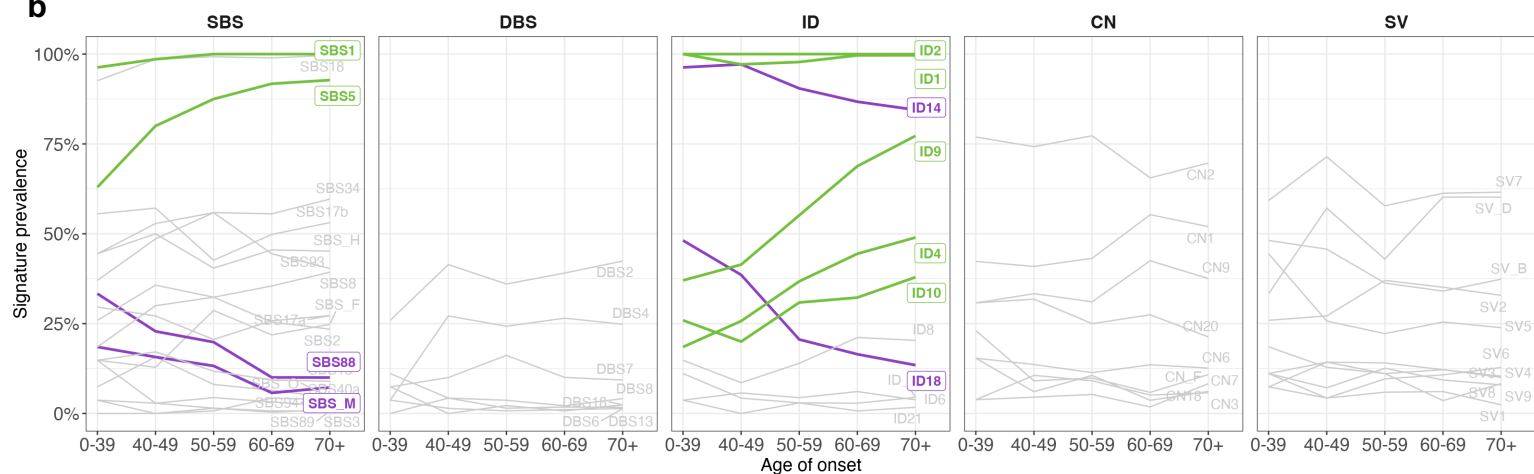
**a**

**Age of onset enrichment** ● Enriched in early-onset patients ● Enriched in late-onset patients

Adjusted by sex, country, tumor subsite, and purity



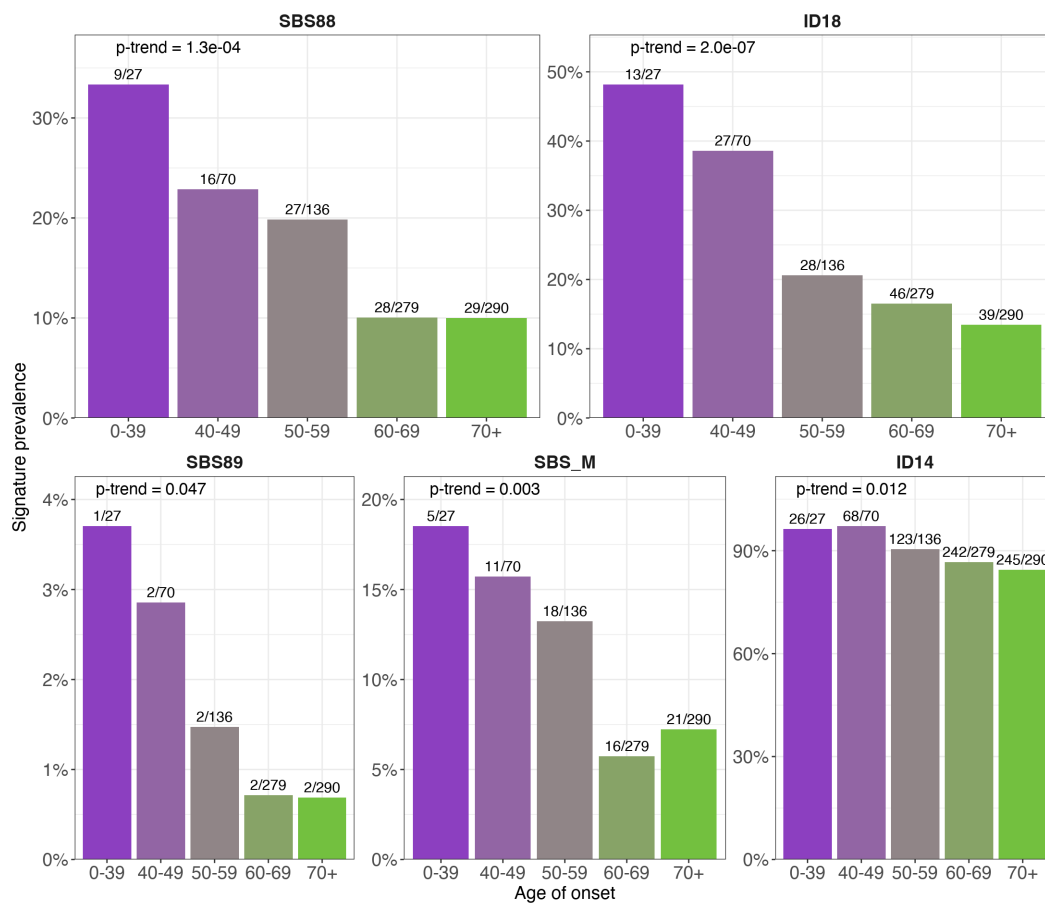
**b**



**c**

**Age of onset trend enrichment**

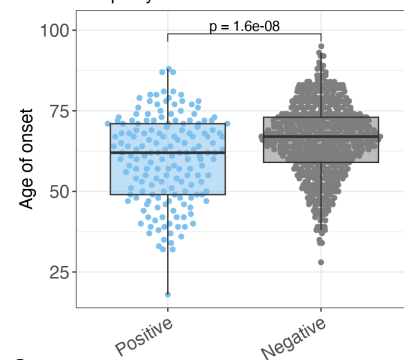
Adjusted by sex, country, tumor subsite, and purity



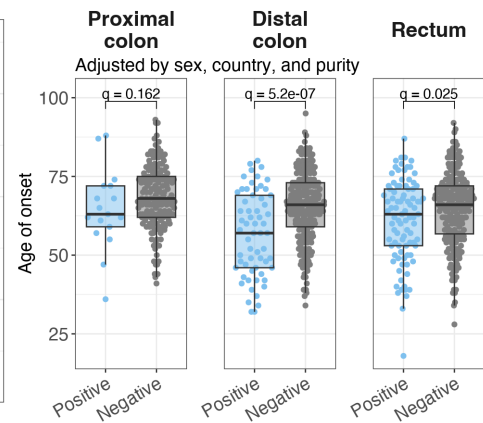
**d**

**Presence of colibactin signatures (SBS88 or ID18)**

Adjusted by sex, country, tumor subsite, and purity

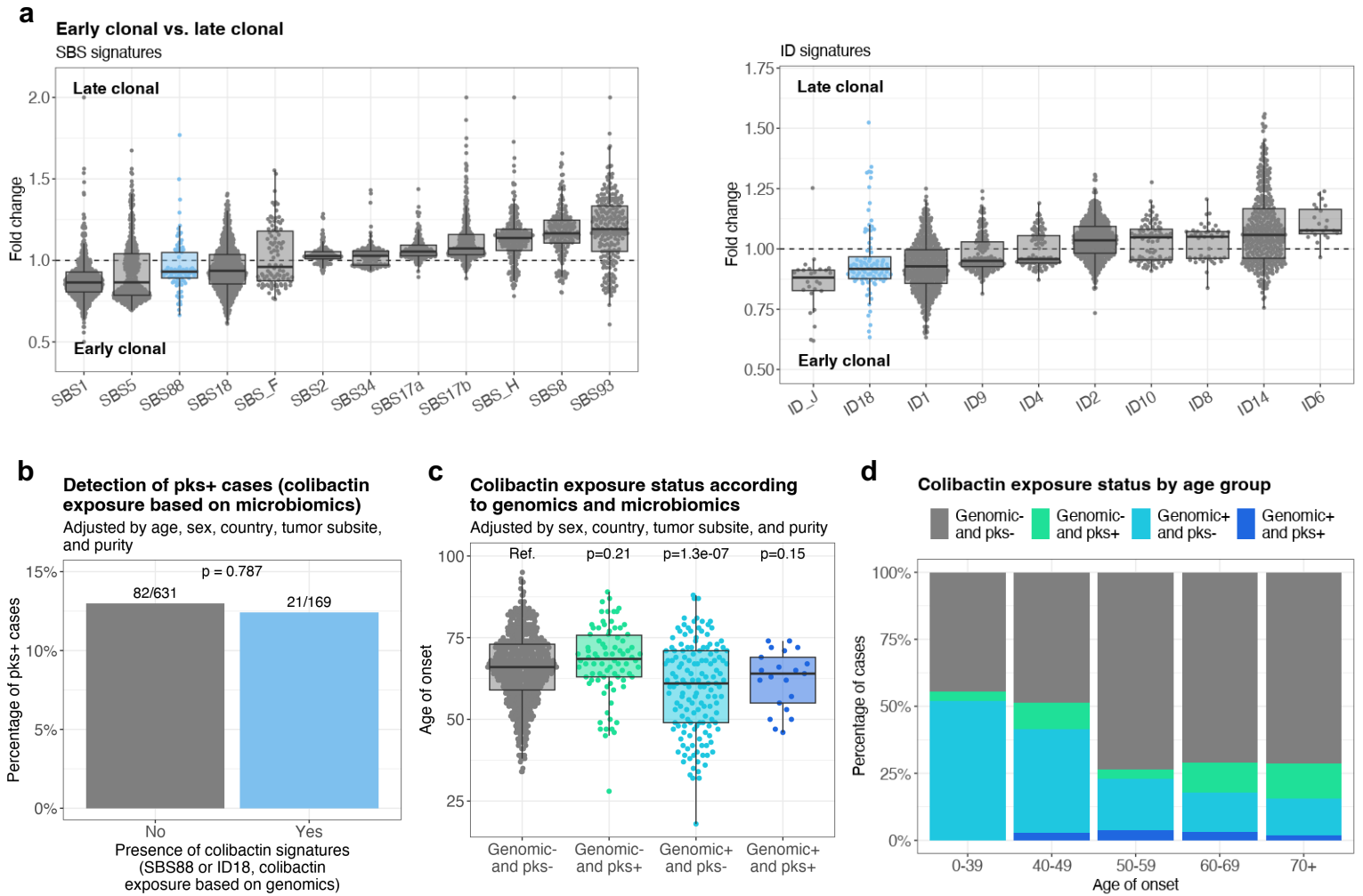


**e**





## Figure 4



## Figure 5

