

The landscape of somatic mutation in normal colorectal epithelial cells

Henry Lee-Six¹, Sigurgeir Olafsson¹, Peter Ellis¹, Robert J. Osborne¹, Mathijs A. Sanders^{1,2},
Luiza Moore¹, Nikitas Georgakopoulos³, Franco Torrente⁴, Ayesha Noorani⁵, Martin
Goddard⁶, Philip Robinson¹, Tim H. H. Coorens¹, Laura O'Neill¹, Christopher Alder¹,
Jingwei Wang¹, Rebecca C. Fitzgerald⁵, Matthias Zilbauer^{4,7}, Nicholas Coleman⁸, Kourosh
Saeb-Parsy³, Inigo Martincorena¹, Peter J. Campbell¹, Michael R. Stratton^{1*}

1. Wellcome Sanger Institute, Hinxton, UK

2. Department of Hematology, Erasmus University Medical Center, Rotterdam, The Netherlands

3. Department of Surgery and Cambridge NIHR Biomedical Research Centre, Cambridge Biomedical Campus, Cambridge, UK

4. Department of Paediatric Gastroenterology, Hepatology, and Nutrition, Addenbrooke's, Cambridge, UK

5. Medical Research Council Cancer Unit, Hutchison/Medical Research Council Research Centre, University of Cambridge, Cambridge, UK, and Cambridge University Hospitals NHS Trust, Hills Road, Cambridge, UK

6. Department of Pathology, Papworth Hospital NHS Trust, UK

7. University Department of Paediatrics, University of Cambridge, UK

8. Department of Pathology, University of Cambridge, Cambridge, UK and Cambridge University Hospitals NHS Foundation Trust, Cambridge, UK

*e-mail: mrs@sanger.ac.uk

Abstract

The colorectal adenoma-carcinoma sequence has provided a paradigmatic framework for understanding the successive somatic genetic changes and consequent clonal expansions leading to cancer. As for most cancer types, however, understanding of the earliest phases of colorectal neoplastic change, which may occur in morphologically normal tissue, is comparatively limited. Here, we whole genome sequenced hundreds of normal crypts from 42 individuals. Signatures of multiple mutational processes were revealed, some ubiquitous and continuous, others only found in some individuals, in some crypts or during certain periods of life. Likely driver mutations were present in ~1% of normal colorectal crypts in middle-aged individuals, indicating that adenomas and carcinomas are rare outcomes of a pervasive process of neoplastic change across morphologically normal colorectal epithelium. Colorectal cancers exhibit substantially elevated mutation burdens relative to normal cells. Sequencing normal colorectal cells provides quantitative insights into the genomic and clonal evolution of cancer.

Introduction

Sequencing of >20,000 cancers has identified the repertoire of driver mutations in cancer genes converting normal cells into cancer cells and revealed the mutational signatures of the underlying biological processes generating somatic mutations^{1,2}. Cancers are, however, end stages of an evolutionary process operating within cell populations and commonly arise through the accumulation of multiple driver mutations engendering a series of clonal expansions. Understanding this progression has depended on identifying somatic mutations in morphologically abnormal neoplastic proliferations representing intermediate stages between normal and cancer cells.³

As for most cancer types, however, the earliest stages of progression to colorectal cancer remain less well understood. The driver mutation that first sets a colorectal epithelial cell on the path to cancer is likely caused by mutational processes operative in normal cells, of which there is limited understanding. The nature and numbers of the earliest neoplastic clones with

53 driver mutations, which conceivably are morphologically indistinguishable from normal
54 cells, are similarly unclear. In large part, these deficiencies are due to the technical challenge
55 of identifying somatic mutations in normal tissues, which are composed of myriad
56 microscopic cell clones. Several different approaches have been adopted to address this⁴⁻¹⁴,
57 revealing signatures of common somatic mutational processes in normal cells of the small
58 and large intestine, liver, blood, skin, and nervous system. Thus far, however, studies have
59 not been of sufficient scale to characterise variation in signature activity or detect less
60 frequent processes⁴⁻¹⁴. Remarkably high proportions of normal skin, oesophageal, and
61 endometrial epithelial cells have been shown to be members of clones already carrying driver
62 mutations^{10,11,15,16}, and large mutant clones have been detected in blood¹⁷⁻²⁰. The extent of
63 this phenomenon in the colon, an organ with a high cancer incidence, has not been
64 investigated.

65
66 Colonic epithelium is a contiguous cell sheet organised into ~15,000,000 crypts each
67 composed of ~2,000 cells²¹. Towards the base of each crypt resides a small number of stem
68 cells ancestral to the maturing and differentiated cells in the crypt²². These stem cells
69 stochastically replace one another through a process of neutral drift^{23,24} such that all stem
70 cells, and thus all cells, in a crypt derive from a single ancestor stem cell that existed in recent
71 years²⁵⁻²⁷. The somatic mutations that were present in this ancestor are thus found in all
72 ~2,000 descendant cells and can be revealed by DNA sequencing of an individual crypt.
73 These stem cells are thought to be the cells of origin of colorectal cancers²⁸. To characterise
74 the earliest stages of colorectal carcinogenesis, somatic mutation burdens, mutational
75 signatures, clonal dynamics, and the frequency of driver mutations in normal colorectal
76 epithelium were explored by sequencing individual colorectal crypts.

77 **Results**

78 **Somatic mutations and mutational signatures**

79 2,035 individual colonic crypts from the normal epithelium of 42 individuals aged 11 to 78,
80 of whom 15 had a history of colorectal cancer and 27 did not (Methods, Supplementary Table
81 1), were isolated using laser capture microdissection and sequenced. The distribution of
82 mutation allele fractions from whole genome sequencing of 571 individual crypts showed
83 that crypts were derived from a single ancestral stem cell (Extended Data Fig. 1d), and
84 simulations indicated that ~90% of mutations called were fully clonal (Supplementary
85 Results 2). There was substantial variation in mutation burdens between individual crypts,
86 ranging from 1,508 to 15,329 for individuals in their sixties, which was not obviously
87 attributable to technical factors. To explore the biological basis of this variation we extracted
88 mutational signatures and estimated the contribution of each to the mutation burden of every
89 crypt (Methods, Supplementary Results 1).

90
91
92 Nine single base substitution (SBS), six doublet base substitution (DBS), and five small indel
93 (ID) mutational signatures were found. Of these, 14 closely matched (Methods) a known
94 reference signature (SBS1, SBS2, SBS5, SBS13, SBS18, DBS2, DBS4, DBS6, DBS8,
95 DBS9, DBS11, ID1, ID2, and ID5, nomenclature as in Alexandrov et al¹) and six did not
96 (SBSA, SBSB, SBSC, SBSD, IDA, and IDB) (Fig. 1, Extended Data Fig. 2-4). Thus, new
97 mutational signatures were extracted despite extensive prior analysis of cancers, perhaps due
98 to masking by the comparative complexity of signature mixtures present in cancer genomes.

99 **Ubiquitous mutational signatures**

100

101 11 signatures (SBS1, SBS5, SBS18, DBS2, DBS4, DBS6, DBS9, DBS11, ID1, ID2, and
102 ID5) were found in >85% of crypts and are here termed “ubiquitous”. All have been
103 previously described¹.

104

105 SBS1 is characterised by C>T substitutions at NCG trinucleotides (the mutated base is
106 underlined) and is likely due to deamination of 5-methylcytosine. Its mutation load correlated
107 linearly with age (Fig. 2). There was, however, variation in SBS1 mutation burdens between
108 crypts from the same individual ($p=2.25e-27$). This was due, in part, to different SBS1
109 mutation rates in different colonic sectors, with mean rates across individuals of 16.8
110 mutations per year (95% CI 15.2-18.3) in the right (ascending and caecum), 16.1 (95% CI
111 14.4-17.5) in the transverse, and 12.8 (95% CI 11.1-14.4) in the left (descending and
112 sigmoid) colon. The SBS1 mutation rate in the terminal ileum was 12.7 (95% CI 10.6-14.9)
113 (Supplementary Results 1). SBS5 is a flat, featureless signature of unknown cause and SBS18
114 is characterised by C>A mutations, which may be due to DNA damage by reactive oxygen
115 species^{29,30}. Their mutation burdens correlated with age, with the same ordering of sector
116 differences as SBS1 ($p=9.89e-26$ for SBS5, $p=5.43e-22$ for SBS18). Even after taking
117 anatomical location and age into account, differences in mutation burden remained between
118 different crypts, notably for SBS18 (Fig. 2, Extended Data Fig. 9, Extended Data Fig. 6al).
119 Combining ubiquitous SBS mutational signatures, and averaging over anatomical sites, the
120 mutation rate was 43.6 mutations per year, comparable with previous estimates⁴.

121

122 DBS2, DBS4, DBS6, DBS9, and DBS11 were tightly correlated in all colonic crypts. ID1,
123 ID2, and ID5, which are characterised by insertions and deletions of a single T and may be
124 the consequence of slippage during DNA replication, all accumulated linearly with age with
125 the same order of sector differences as SBS1 ($p=1.66e-05$ for ID1, $p=4.53e-06$ for ID2, and
126 $p=4.53e-06$ for ID5) (Supplementary Results 1, Extended Data Fig. 5).

127

128 The correlations of ubiquitous signatures with age indicate that the mutational processes
129 underlying them operate throughout life, in all individuals and all colorectal stem cells.
130 However, the results also suggest that differences in physiology and/or microenvironment
131 (and potentially age of the most recent common ancestor of crypts²⁷) between different
132 sectors of the colon cause measurable differences in somatic mutation rates.

133

134 **Sporadic mutational signatures**

135 Nine signatures (SBS2, SBS13, SBSA-D, DBS8, and IDA-B) were present only in a subset
136 of individuals and/or a subset of crypts and are termed “sporadic”. All were novel, except for
137 SBS2, SBS13 and DBS8. SBS2 and SBS13 are characterised by C>T and C>G mutations at
138 TCN, are likely due to APOBEC cytidine deaminases and usually occur together^{31,32}. They
139 were unequivocally present in only two crypts (a colonic crypt (Extended Data Fig. ai) and an
140 ileal crypt (Extended Data Fig. ao) from different individuals), occurring together and each
141 accounting for over 150 mutations. To our knowledge, this is the first report that APOBEC
142 DNA-editing of the human genome occurs in normal cells *in vivo*. The sequence context of
143 these mutations in normal colon suggests that APOBEC3A is the major contributing
144 enzyme³³.

145

146 Four SBS signatures that do not match the reference set, SBSA-D, were found in normal
147 colorectal cells (SBSA has recently been reported in an oral squamous carcinoma³⁴). SBSA is
148 characterised by T>C at ATA, ATT, and TTT, and T>G at TTT. Its mutation burden
149 correlated closely with that of IDA, in which single T deletions in short runs of Ts (with a
150 mode of four) predominate, suggesting that they are due to the same underlying mutational

151 process. SBSA was detectable in 29/42 individuals, often accounting for thousands of
152 mutations in just a subset of crypts. It clustered spatially in the colon, with crypts from the
153 same biopsy carrying the signature even though the mutations themselves were not shared
154 (Supplementary Results 1, Extended Data Fig. 9). 2.5-fold more T>C mutations occurred
155 when the T was on the transcribed than on the untranscribed strand. Transcriptional strand
156 bias is often due to transcription coupled nucleotide excision repair acting on DNA damaged
157 by exogenous exposures causing covalently bound bulky adducts, but can also be caused by
158 transcription coupled DNA damage³⁵. Assuming either is the case, damage to adenine
159 underlies SBSA. To investigate the timing of SBSA, phylogenetic trees of mutations were
160 constructed and the mutational signatures in each branch established (Fig. 3, Extended Data
161 Fig. 6). SBSA was confined to early branches of these phylogenies (when these were
162 available for analysis) (Fig. 3b, Extended Data Fig. 6 f, h, z, aa, am, ao, aq). Using the
163 number of SBS1 mutations as indicators of real time, the mutational process underlying
164 SBSA appears to be active before 10 years of age (Supplementary Methods, Extended Data
165 Fig. 6aq). SBSA may therefore be caused by an extrinsic, locally acting and patchily
166 distributed mutagenic insult occurring during childhood.

167
168 SBSB was characterised by C>T at ACA, T>A at CTN, and T>G at GTG and was present in
169 subsets of crypts from four individuals (e, aa, ai, and aj in Extended Data Figure 6),
170 accounting for variable numbers of substitutions, with a maximum of 3,002 in one crypt.
171 (Fig. 3c, Extended Data Fig. 6ai). In the two individuals in whom it could be timed (Extended
172 Data Fig. 6 aa, ai, aq), it appeared – as with SBSA – to be most active in the first decade of
173 life. SBSB correlated with DBS8 and IDB (Fig. 3c, Extended Data Fig. 9), suggesting that
174 they are caused by the same underlying mutational process. DBS8 is composed of AC>CA
175 and AC>CT mutations and has previously been reported in rare hypermutated cancers with
176 no obvious cause¹. IDB is dominated by deletion of a single T with no other Ts surrounding
177 it.

178
179 SBSC is characterised by one C>T mutation in CC dinucleotides. It primarily affects three
180 crypts, with 1,050, 827, and 695 mutations respectively, from the left colon of one individual
181 with an unremarkable history (Extended Data Fig. 9, Extended Data Fig. 6m, Supplementary
182 Table 1).

183
184 All crypts from a 66 year-old man carried many thousands of mutations of SBSD (Figure 3d,
185 Extended Data Fig. 6ap), characterised by T>A substitutions with a transcriptional strand bias
186 compatible with damage to adenine. This individual had been treated with multiple
187 chemotherapeutic agents (cyclophosphamide, doxorubicin, vincristine, prednisolone,
188 chlorambucil, bleomycin and etoposide) for lymphoma and subsequently developed caecal
189 adenocarcinoma. SBSD resembles SBS25 (cosine similarity 0.9), previously found in
190 Hodgkin lymphoma cell lines from two chemotherapy-treated patients^{31,36}. To our knowledge
191 this is the first time that the mutational consequences of chemotherapy have been
192 demonstrated in normal human cells *in vivo*. The mutation burden in his colorectal epithelium
193 was 3-5 fold higher than expected for his age, thus by extrapolation equivalent to that of a
194 200-300 year-old.

195 196 **Copy number changes and structural variants**

197 Copy number changes and/or structural variants were found in 80 out of 449 (18%) evaluable
198 crypts. Five crypts exhibited eight whole chromosome copy number increases which affected
199 the same three chromosomes – 3, 7 and 9 – as well as the X chromosome (Extended Data
200 Fig. 7a). Thus, copy number increases clustered in certain crypts and tended to affect certain

201 chromosomes. No whole chromosome losses were observed. Arm-level chromosome 7 copy
202 number increases are common in colorectal cancers³⁷ and adenomas³⁸. Chromosome 3 and 9
203 copy number increases are seen in colorectal cancers, but are almost as frequently deleted³⁷.
204 Copy number neutral loss of heterozygosity (CNN-LOH) was observed in 12 crypts,
205 affecting chromosomes 1p, 6p, 7p, 8q, 9q, 10q (twice), 17p, 17q, 18q, 21q and 22q (Extended
206 Data Fig. 7c). CNN-LOH is frequently observed in colorectal cancers, although the specific
207 changes that we observe here are not recurrent features³⁹. Five copy number changes could be
208 timed and all were estimated to have occurred in adulthood (Extended Data Fig. 7b). Two
209 changes that affected the same crypt appeared to be synchronous (Supplementary Results 1).
210 Structural variant analysis detected 48 large deletions, 18 tandem duplications, four
211 translocations, and two inversions (Extended Data Fig 7d, Supplementary Results 1). All
212 were private to a single crypt, except for one deletion which was present in two adjacent
213 crypts sharing few mutations, indicating that it occurred during gestation or early childhood.

214 215 **Driver mutations**

216 Driver mutations are those that confer a selective advantage during cancer evolution. To
217 search for drivers in normal colon, the whole genome sequences of 571 crypts were
218 supplemented with targeted sequencing of 90 known colorectal cancer genes (Supplementary
219 Table 4) in additional crypts. In total, substitutions in these genes were evaluable in 1,403
220 crypts and indels in 1,046. Statistical analysis revealed evidence of positive selection on the
221 recessive cancer genes *AXIN2* (three truncating mutations, adjusted q value 0.004) and
222 *STAG2* (two truncating mutations, adjusted q value 0.038) indicating that these mutations are
223 likely drivers. Additional likely drivers were identified in cancer genes with canonical
224 missense hotspot mutations. Nine hotspot mutations in *PIK3CA* (E542K, R38H), *ERBB2*
225 (R678Q, V842I, T862A), *ERBB3* (R475W, R667L), and *FBXW7* (R505C, R658Q) were
226 observed (Extended Data Fig. 8). Given the specificity of these hotspot mutations, most are
227 likely to be drivers. In addition, heterozygous truncating mutations were found in the
228 recessive cancer genes *ARID2*, *ATM* (two), *ATR*, *BRCA2*, *CDK12* (two), *CDKN1B*, *RNF43*
229 (two), *TBLIXR1*, and *TP53* (Supplementary Table 5). There was no statistical evidence for
230 selection of truncating mutations in the set of 90 colorectal cancer genes overall. The
231 possibility that some have conferred clonal growth advantage, however, is not excluded.
232 None of the analysed crypts carried more than one putative driver.

233
234 23 pairs of adjacent crypts shared over 100 SBS1 mutations and thus were likely to have been
235 generated by postnatal crypt fission. Two pairs carried driver mutations (one *AXIN2* nonsense
236 mutation and one *PIK3CA* E542K), although the association of driver mutations with crypt
237 fission is not significant (p=0.17). In one sister crypt the *AXIN2* mutation was rendered
238 homozygous by CNN-LOH of 17q, revealing ongoing clonal evolution in normal colon (Fig.
239 4, Fig. 3b).

240
241 On the conservative assumption that just the *AXIN2* and *STAG2* truncating mutations and the
242 missense hotspot mutations in *PIK3CA*, *ERBB2*, *ERBB3* and *FBXW7* are drivers, ~1% of
243 normal colorectal crypts (~150,000 crypts) in a 50-60 year old (the mean age of crypts
244 assessed for drivers in our cohort was 53 years old) carries a driver mutation. Since in the
245 over 70s ~40% of people have an adenoma on colonoscopy⁴⁰ and ~5% of people develop
246 colorectal cancer over their lifetime⁴¹ (and some of these may arise from more recently-
247 acquired driver mutations) only an extremely small proportion of these crypt microneoplasms
248 becomes a macroscopically detectable adenoma (< 1/375,000) or carcinoma (< 1/3,000,000)
249 within the following few decades.

250

251 **Clonal dynamics of normal epithelium**

252 The distribution of allele fractions of mutations within the crypt informs on the dynamics of
253 stem cell turnover within the crypt. We estimate that the average number of years to the most
254 recent common ancestor of crypts is 5.5 years (CI95 for the mean: 1-10.5 years), similar to
255 previous estimates¹⁴. Our data are compatible with previous estimates of 7 active stem cells
256 and 1.3 stem cell replacements per year¹⁴ or with 5 stem cells and 0.6 stem cell replacements
257 per year⁴², but we cannot exclude a larger number of stem cells turning over more frequently
258 (Extended data figure 10, Supplementary Results 2). The microdissection approach also
259 allowed investigation of the clonal structure of colonic epithelium beyond the crypt. By
260 comparing the genetic relatedness of crypts with their spatial relatedness, we estimate that
261 crypts fission at a mean rate of once every 27 years (CI95: 15.9-47.6 years) (Extended data
262 figure 10, Supplementary Results 2).

263

264 **Comparisons with colorectal cancer**

265 There are marked differences between the genomes of normal colorectal stem cells and those
266 of colorectal cancers. The total mutation burdens of substitutions (10,000-20,000) and indels
267 (1,000-2,000) found in most colorectal carcinomas¹ (excluding those with hypermutator
268 phenotypes in which it is usually >10-fold more) are higher than the ~3,000 substitutions and
269 300 indels found in most normal crypts from 50-60 year-old individuals (Extended Data
270 Figure 11a). These differences may be underestimated as the most recent common ancestor
271 of cancers likely predates that of normal crypts. The high mutation burdens and associated
272 mutational signatures of DNA mismatch repair deficiency and/or polymerase ϵ/δ mutations
273 were not found in any normal colorectal crypts but are present in ~20% colorectal cancers.
274 Equally striking is the difference between the 0-4 structural changes per normal crypt (with
275 the majority having none (Supplementary Results 1)) and the 10s to 100s per colorectal
276 cancer⁴³. In all these respects, the genomes of normal crypts with driver mutations were
277 similar to those of normal crypts without drivers (Extended Data Fig. 9).

278

279 There was no difference in the burden either of sporadic or ubiquitous mutational processes
280 between the crypts of individuals with and without a colorectal cancer (Supplementary
281 Results 1). If differences in mutational processes in normal cells do underlie why some
282 people develop colon cancer and others do not, these mutational processes must affect only a
283 small proportion of crypts in the colon, or only exert subtle effects on the mutation rate such
284 that we could not detect differences between the two groups. The increased base substitution
285 and indel mutation loads in cancers are due to a combination of higher burdens of the
286 ubiquitous mutational signatures found in normal crypts, additional signatures thus far found
287 exclusively in cancers (confirming previous reports^{5,44}) and larger numbers of copy number
288 changes and structural variation (Extended Data Figure 11a). The causes of some of these
289 additional mutations in cancer are known (for example, defective mismatch repair and
290 polymerase ϵ/δ mutations) but the majority are uncertain.

291

292 The relative frequencies of mutated cancer genes differ between colorectal
293 adenomas/carcinomas and normal colorectal cells ($p=0.003$, Supplementary Results 1,
294 Extended Data Figure 11a). In colorectal cancer, mutations in *APC*, *KRAS* and *TP53* are
295 common³⁷, accounting for 56% of base substitution and indel drivers (Supplementary
296 Methods) but are comparatively rare among normal crypts with driver mutations (1/14). By
297 contrast, mutations in, for example, *ERBB2* and *ERBB3* are common in normal crypts with
298 drivers (5/14) but rare in colorectal cancer (7/631). In the case of *APC* (but not *KRAS* and
299 perhaps not *TP53*), biallelic inactivation may be required to confer a strong growth
300 advantage, which helps to explain why *APC* may be mutated less frequently in normal colon

301 than *ERBB2/3* that require a single hit to do so. The results suggest that mutations in *APC*,
302 *KRAS* and *TP53* confer higher likelihoods of conversion to adenoma and carcinoma than
303 mutations in *ERBB2* and *ERBB3* whereas the latter confer higher likelihoods of stem cells
304 colonising crypts. There was no detectable difference in the frequency of driver mutations
305 between individuals in our cohort who had colorectal cancer relative to those who did not
306 (Supplementary Results 1).

307
308

309 **Discussion**

310

311 This study has characterised all classes of somatic mutation in hundreds of normal colorectal
312 epithelial stem cells. Our experimental design allows us to gain insights into different facets
313 of the earliest stages of the clonal evolution of colorectal cancers, namely the range of
314 mutational processes, the frequency of driver mutations, and the clonal dynamics of colonic
315 stem cells.

316

317 A substantial repertoire of base substitution and indel mutational processes is operative, some
318 ubiquitous and some sporadic, together with relatively infrequent copy number changes and
319 genome rearrangements. APOBEC DNA-editing occurs in normal colon, albeit only in rare
320 cells. Many signatures, however, are of unknown aetiology and some appear to be acquired
321 early in life. The presence of five times the age-standard mutation load in all colorectal cells,
322 and potentially many other tissues, in an individual who had undergone chemotherapy
323 provides new insight into the impact of such exposures and raises questions pertaining to its
324 relationship with chemotherapy's relatively modest impact on cancer risk⁴⁵.

325

326 The earliest stages of colorectal cancer development have been revealed in this manuscript.
327 They are characterised by numerous crypts carrying driver mutations, of which only a very
328 small fraction ever manifests as macroscopic neoplasms. Certain mutated cancer genes
329 appear to foster this pervasive and invisible wave of microneoplastic change whereas others
330 particularly engender progression to colorectal adenoma and cancer. The conversion of these
331 early microneoplasms to more advanced stages of colorectal neoplasia is associated with
332 acquisition of elevated mutational loads, composed of base substitutions, indels, structural
333 variants and copy number changes. More extensive studies of colorectal epithelium will
334 enable characterisation of the rarer intermediate stages between these early clones and small
335 adenomas, and refine understanding of the development of the subset of microneoplasms
336 with higher likelihoods of becoming carcinomas.

337

338 The proportion of normal colorectal epithelial cells with driver mutations (1%) is, however,
339 substantially lower than that of other normal tissues so far studied, notably skin (30%)¹⁰,
340 oesophagus (>50%)¹⁶. This may be due, at least in part, to the modular structure of glandular
341 epithelia. The small number of stem cells within a crypt diminishes the probability that a cell
342 with a driver mutation will outcompete its wild-type neighbours. Moreover, even if it does
343 colonise the crypt, a mutant stem cell is entombed in it unless it can overcome the largely
344 unknown forces that govern clonal expansion through crypt fission. The lower driver burden
345 of colon relative to endometrium^{11,15}, which is also glandular, remains to be explored.

346

347 Fundamental questions are being addressed with respect to differences in cancer incidence
348 rates between tissues. The somatic mutation burden in colon and ileum is similar despite the
349 substantially higher cancer incidence rate in colon (as previously noted⁴) and therefore does
350 not appear to account for this difference. Whether the total burden of microneoplastic change

351 across the colon and in other tissues more closely correlates with these differences is yet to be
352 determined.

353

354 Finally, this study provides a reference perspective on the mutational signatures and driver
355 mutations in normal colon against which disease states of inflammatory, genetic, neoplastic,
356 degenerative and other aetiologies can be compared. Similar surveys conducted across the
357 range of normal cell types will inform on the universal process of somatic evolution in the
358 human body in health and disease.

359

360

361

362 **ACKNOWLEDGEMENTS**

363 This work was supported by the Wellcome Trust. We thank Paul Scott, Jo Fowler, David
364 Fernandez-Antoran, and Yvette Hooks for their advice with histology and laser capture
365 microdissection, and Moritz Gerstung for his advice on statistics. We thank the Sanger
366 Institute Research and Development Facility for their help sequencing microbiopsies. The
367 authors would like to thank the staff of WTSI Sample Logistics, Genotyping, Pulldown,
368 Sequencing and Informatics facilities for their contribution. We further thank Krishnaa
369 Mahbubani, Rogier ten Hoopen, Cinzia Scarpini, and the Phoenix study team of Nicola
370 Grehan, Irene Debiram-Beecham, Jason Crawte, Tara Nukcheddy Grant, Pierre Lao-Sirieix,
371 and Andy Hindmarsh for their help with sample collection. Access to transplant organ donor
372 samples was provided by the Cambridge Biorepository for Translational Medicine. Ayesha
373 Noorani was funded through an MRC Clinical Research Fellowship. The autopsy cohort
374 was funded through this and an MRC core grant (RG84369) and an NIHR
375 Research Professorship (RG67258) to Rebecca Fitzgerald. We thank the Human Research
376 Tissue Bank, which is supported by the National Institute for Health Research (NIHR)
377 Cambridge Biomedical Research Centre, from Addenbrooke's Hospital. Additional
378 infrastructure support was provided from the CRUK funded Experimental Cancer Medicine
379 Centre in Cambridge. Finally, we thank all the individuals who contributed samples to this
380 study.

381

382

383 **AUTHOR CONTRIBUTIONS**

384 MRS and HLS designed the study and wrote the manuscript with contributions from all the
385 authors. KSP, NC, MZ, RCF, NG, FT, AN, MG, and LM recruited patients and obtained
386 samples. PE, RO, HLS, and LM devised the protocol to laser capture microdissect and
387 sequence colonic crypts. HLS prepared sections, microdissected, and lysed colonic crypts. PR
388 contributed to laser capture microdissection. PE and CA made libraries. HLS performed most
389 of the data curation and statistical analysis. SO estimated the rate of crypt fission. MAS
390 devised filters for substitution calling. JW performed in-house NMF signature extraction. TC
391 and PR contributed to statistical analyses. LON provided technical assistance. PJC and IM
392 oversaw statistical analyses. MRS supervised the study.

393

394 **COMPETING INTERESTS**

395 The authors declare no competing interests.

396

397

398 **REFERENCES**

399

- 400 1. Alexandrov, L.B. et al. The repertoire of mutational signatures in human cancer.
401 Preprint at: <https://www.biorxiv.org/content/early/2018/05/15/322859> (2018).
- 402 2. Sabarinathan, R. et al. The whole genome panorama of cancer drivers. Preprint at:
403 <https://www.biorxiv.org/content/early/2017/12/23/190330> (2017).
- 404 3. Fearon E.R. & Vogelstein B. A genetic model for colorectal tumorigenesis. *Cell* **61**,
405 759-767 (1990).
- 406 4. Blokzijl, F. et al. Tissue-specific mutation accumulation in human adult stem cells
407 during life. *Nature* **538**, 260–264 (2016).
- 408 5. Roerink S.F. et al., Intra-tumour diversification in colorectal cancer at the single cell
409 level. *Nature* **556**, 457-462 (2018).
- 410 6. Welch, J. S. et al. The origin and evolution of mutations in acute myeloid leukemia.
411 *Cell* **150**, 264–278 (2012).
- 412 7. Bae, T. et al. Different mutational rates and mechanisms in human cells at
413 pregastrulation and neurogenesis. *Science* **359**, 550–555 (2018).
- 414 8. Behjati, S. et al. Genome sequencing of normal cells reveals developmental lineages
415 and mutational processes. *Nature* **513**, 422–425 (2014).
- 416 9. Lee-Six H. et al., Population dynamics of normal human blood inferred from somatic
417 mutations. *Nature* <https://doi.org/10.1038/s41586-018-0497-0> (2018).
- 418 10. Martincorena, I. et al. Tumor evolution. High burden and pervasive positive
419 selection of somatic mutations in normal human skin. *Science* **348**, 880–886
420 (2015).
- 421 11. Suda, K. et al. Clonal Expansion and Diversification of Cancer-Associated
422 Mutations in Endometriosis and Normal Endometrium. *Cell Rep.* **24**, 1777-1789
423 (2018).
- 424 12. Lodato, M. A. et al. Aging and neurodegeneration are associated with increased
425 mutations in single human neurons. *Science* **359**, 555–559 (2018).
- 426 13. Hoang, M.L. et al. Genome-wide quantification of rare somatic mutations in normal
427 human tissues using massively parallel sequencing. *PNAS* **113**, 9846-9851 (2016).
- 428 14. Nicholson, A. et al. Fixation and spread of somatic mutations in adult human colonic
429 epithelium. *Cell Stem Cell* **22**, 909-918 (2018).
- 430 15. Moore, L. et al. The mutational landscape of normal human endometrial epithelium.
431 Preprint at: [https://www.biorxiv.org/content/10.1101/505685v1\(2018\)](https://www.biorxiv.org/content/10.1101/505685v1(2018)).
- 432 16. Martincorena, I. Somatic mutant clones colonize the human esophagus with age.
433 *Science* **362**, 911-917 (2018).
- 434 17. Jaiswal, S. et al. Age-related clonal hematopoiesis associated with adverse outcomes.
435 *N. Engl. J. Med.* **371**, 2488–2498 (2014).
- 436 18. Xie, M. et al. Age-related mutations associated with clonal hematopoietic expansion
437 and malignancies. *Nat. Med.* **20**, 1472–1478 (2014).
- 438 19. McKerrell, T. et al. Leukemia-associated somatic mutations drive distinct patterns of
439 age-related clonal hemopoiesis. *Cell Rep.* **10**, 1239–1245 (2015).
- 440 20. Genovese, G. et al. Clonal hematopoiesis and blood-cancer risk inferred from blood
441 DNA sequence. *N. Engl. J. Med.* **371**, 2477-2487 (2015).
- 442 21. Potten, C.S. et al. Measurement of in vivo proliferation in human colorectal mucosa
443 using bromodeoxyuridine. *Gut.* **33**, 71-78 (1992).
- 444 22. Cheng, H. & Leblond, C.P. Origin, differentiation and renewal of the four main
445 epithelial cell types in the mouse small intestine. V. Unitarian Theory of the origin of
446 the four epithelial cell types. *Am J Anat.* **141**, 537-561 (1974).
- 447 23. Lopez-Garcia, C. et al. Intestinal stem cell replacement follows a pattern of neutral
448 drift. *Science* **330**, 822-825 (2010).

- 449 24. Snippert, H.J. et al. Intestinal crypt homeostasis results from neutral competition
450 between symmetrically dividing Lgr5 stem cells. *Cell* **143**, 134-144 (2010).
- 451 25. Griffiths D.F. et al. Demonstration of somatic mutation and colonic crypt clonality by
452 X-linked enzyme histochemistry. *Nature* **333**, 461-463 (1988).
- 453 26. Winton, D.J., and Ponder, B.A. Stem-cell organization in mouse small intestine. *Proc.*
454 *Biol. Sci.* **241**, 13-18 (1990).
- 455 27. Kozar, S. et al. Continuous clonal labeling reveals small numbers of functional stem
456 cells in intestinal crypts and adenomas. *Cell Stem Cell* **13**, 626-633 (2013).
- 457 28. Barker, N. et al. Crypt stem cells as the cells-of-origin of intestinal cancer. *Nature*
458 **457**, 608-611 (2009).
- 459 29. Rouhani, F.J. et al. Mutational history of a human cell lineage from somatic to
460 induced pluripotent stem cells. *PLoS Genet.* **12**, e1005932 (2016).
- 461 30. Viel, A. et al. A specific mutational signature associated with DNA 8-Oxoguanine
462 persistence in MUTYH-defective colorectal cancer. *EBioMedicine* **20**, 39-49 (2017).
- 463 31. Alexandrov, L.B. et al. Signatures of mutational processes in human cancer. *Nature*
464 **500**, 415-421 (2013).
- 465 32. Nik-Zainal, S. et al. Mutational processes molding the genomes of 21 breast cancers.
466 *Cell* **149**, 979-993 (2012).
- 467 33. Chan, K. et al. An APOBEC3A hypermutation signature is distinguishable from the
468 signature of background mutagenesis by APOBEC3B in human cancers. *Nat Genet.*
469 **47**, 1067-1072 (2015).
- 470 34. Boot, A. et al. Mutational signature analysis of Asian OSCCs reveals novel
471 mutational signature with exceptional sequence context specificity. Preprint at:
472 <https://www.biorxiv.org/content/early/2018/07/19/368753.1> (2018).
- 473 35. Haradhvala, N. J. et al. Mutational Strand Asymmetries in Cancer Genomes Reveal
474 Mechanisms of DNA Damage and Repair. *Cell* **164**, 538-549 (2016).
- 475 36. Wolf, J. et al. Peripheral blood mononuclear cells of a patient with advanced
476 Hodgkin's lymphoma give rise to permanently growing Hodgkin-Reed Sternberg
477 cells. *Blood* **87**, 3418-3428 (1996).
- 478 37. The Cancer Genome Atlas Network, Comprehensive molecular characterization of
479 human colon and rectal cancer. *Nature* **487**, 330-337 (2012)
- 480 38. Bomme, L. et al. Cytogenetic analysis of colorectal adenomas: karyotypic
481 comparisons of synchronous tumors. *Cancer Genet Cytogenet* **106(1)**:66-71 (1998).
- 482 39. Andersen, C.L. et al. Frequent occurrence of uniparental disomy in colorectal cancer.
483 *Carcinogenesis* **28(1)**: 38-48 (2007).
- 484 40. Corley, D.A., et al. Variation of adenoma prevalence by age, sex, race, and colon
485 location in a large population: implications for screening and quality programs. *Clin.*
486 *Gastroenterol. Hepatol.* **11**, 172-180 (2013).
- 487 41. Cancer Research UK, Bowel Cancer Incidence Statistics,
488 [https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-](https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/bowel-cancer/incidence#heading-Seven)
489 [cancer-type/bowel-cancer/incidence#heading-Seven](https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/bowel-cancer/incidence#heading-Seven) (Accessed August 2018).
- 490 42. Stamp, C. et al. Predominant asymmetrical stem cell fate outcome limits the rate of
491 niche succession in human colonic crypts. *EBioMedicine* **31**: 166-73 (2018).
- 492 43. Li, Y. et al. Patterns of structural variation in human cancer. Preprint at:
493 <https://www.biorxiv.org/content/early/2017/08/27/181339> (2017).
- 494 44. Lugli, N. et al. Enhanced Rate of Acquisition of Point Mutations in Mouse Intestinal
495 Adenomas Compared to Normal Tissue. *Cell Reports* **19**, 2185-2192 (2017).
- 496 45. Travis, L.B. Therapy-associated solid tumors, *Acta Oncologica*, **41**, 323-333 (2002).
- 497
- 498

499 **FIGURE LEGENDS**

500

501 **Figure 1. Mutational signatures present in normal colon. a**, an example SBS, DBS, and
502 ID signature showing the categories into which mutations are divided. Later figures are
503 shown in the same format. **b**, the complement of signatures in normal colonic epithelium.
504 Known signatures are labelled according to their nomenclature in PCAWG, while novel
505 signatures are labelled with letters. SBS, single base substitution; DBS, doublet base
506 substitution; ID, small insertion or deletion.

507

508 **Figure 2. Mutation burden versus age for every signature.** For every signature, the median
509 (horizontal bar) and range (vertical bar) in mutation burden for all the crypts from each
510 individual are shown. Each individual is coloured differently. n=445 crypts from 42
511 individuals.

512

513 **Figure 3. Crypt phylogenies.** For selected individuals (**a-d**), each phylogeny is shown three
514 times, with branch lengths proportional to SBS (top), DBS (middle), and ID (bottom)
515 mutation counts. A stacked barplot of the signatures contributing to each branch is
516 superimposed. The ordering of signatures along branches is for visualisation purposes. “X0”
517 indicates mutations that could not confidently be assigned to any signature. Phylogenies for
518 all individuals are shown in Extended Data Fig. 6. Selected phylogenies are dominated by
519 ubiquitous signatures (**a**), SBSA and IDA (**b**), SBSB, DBS8, and IDB (**c**), or SBSB for the
520 individual exposed to chemotherapy (**d**).

521

522 **Figure 4. An AXIN2 inactivating mutation. (a)** a section after dissection. Red dots
523 represent crypts with the *AXIN2* mutation, blue dots those without it. Crypts without dots
524 failed sequencing. **(b)** crypts with the mutation appeared no different to others. **(c)** CNN-
525 LOH of one crypt over the *AXIN2* locus. The copy number state (y axis) for every
526 chromosome is shown, with one allele coloured red and the other green. **(d)** Jbrowse image of
527 reads supporting the *AXIN2* mutations. The mutation is red. 25/29 reads support it in the crypt
528 with CNN-LOH; the four that do not presumably represent stromal contamination.

529

530

531 **EXTENDED FIGURE LEGENDS**

532

533 **Extended Data Figure 1. Laser capture microdissection of crypts. (a)** a representative
534 image of a section of colonic tissue, with a magnified inset showing the section before and
535 after dissection of a crypt. **(b-c)**, the coverage of crypts that underwent whole genome **(b)** and
536 targeted **(c)** sequencing. **(d-e)**, their respective VAF (which is half of the clonal fraction). **(f-**
537 **g)**, substitutions **(f)** and indels **(g)** removed by filtering steps and their mutational spectrum,
538 arranged as in Figure 1.

539

540 **Extended Data Figure 2. HDP signature extraction results.** Results of signature
541 extraction using an HDP with conditioning on signatures known to be active in colorectal
542 cancer. For each signature, the extracted signature and the profile of a sample that has a
543 strong contribution of that signature are shown. Signatures are presented as in Fig. 2. The
544 HDP extraction was followed by deconvolution by Expectation Maximisation (Methods,
545 Extended Data Fig. 3) to produce the version of signatures presented in the main text. HDP,
546 Hierarchical Dirichlet Process.

547

548 **Extended Data Figure 3. Expectation maximisation decomposition of HDP signatures.**
549 Three signatures were decomposed. For each panel, the original HDP version is shown on
550 the top left, the PCAWG signatures that are deemed to contribute at least 10% of mutations to
551 it on the right, and the reconstituted signature built by combining the PCAWG signatures on
552 the bottom left. The cosine similarity of the reconstituted signature to the original is shown in
553 the title to the reconstituted signature plot. HDP, Hierarchical Dirichlet Process; PCAWG,
554 Pan Cancer Analysis of Whole Genomes.

555
556 **Extended Data Figure 4. Validation of single base substitution signatures.** Other methods
557 of signature extraction were run to test the robustness of signature decomposition. **a**, HDP
558 without pre-conditioning on PCAWG. **b**, In-house NNMF without pre-conditioning on
559 PCAWG. **c**, NNMF implemented by the MutationalPatterns R package (Methods). HDP,
560 Hierarchical Dirichlet Process; PCAWG, Pan Cancer Analysis of Whole Genomes; NNMF,
561 Non-Negative Matrix Factorisation.

562
563 **Extended Data Figure 5. Linear modelling of signature accumulation.** For signatures that
564 appeared to show a linear accumulation with age, the mutation rate per site was determined
565 using mixed models, with age and site as fixed effects, and individual as a random effect.
566 Confidence intervals were determined by bootstrapping. n=445 crypts from 42 individuals.
567 Solid lines represent the mean slope of the regression and shaded areas its 95% confidence
568 intervals.

569
570 **Extended Data Figure 6. Crypt phylogenies. (a-ap)** For every individual, the phylogeny of
571 crypts is shown three times: on top, with branch lengths proportional to the number of single
572 base substitutions; in the middle, with branch lengths proportional to the number of doublet
573 base substitutions; on the bottom, with branch lengths proportional to the number of small
574 insertions and deletions. Scale bars are shown on the right-hand side. A stacked barplot of the
575 mutational signatures that contribute to each branch is overlaid over every branch. “X0”
576 indicates mutations that could not confidently be assigned to any signature. Please note that
577 the ordering of signatures along a given branch is just for visualisation purposes: we cannot
578 distinguish the timing of different signatures along a branch. **(aq)**, Cumulative burden of
579 SBSA (top panel) and SBSB (bottom panel) is plotted relative to the cumulative burden of
580 SBS1 in order to time these mutational processes throughout life. Informative clades are
581 shown (from patients with labelling as in the other panels), with every node and tip of the
582 clade plotted in the space of the cumulative number of mutations due to a given signature that
583 have occurred up until that node in the tree. Lines represent the branching structure of the
584 tree.

585
586 **Extended Data Figure 7. Copy number changes and structural variants in normal**
587 **colon.** 449 crypts had sufficient coverage to be evaluated. **(a)** whole chromosome
588 amplifications in five crypts. The copy number state (y axis) for each allele, one coloured red,
589 and one coloured green, is shown. Chromosomes are labelled along the top of the graph. **(b)**
590 timing of copy number changes throughout life. Vertical bars represent 95% confidence
591 intervals determined by bootstrapping. Horizontal bars represent the most likely time of the
592 copy number change, as defined by mutationTimeR (Supplementary Results 1) **(c)**, crypts
593 with loss of heterozygosity. For each chromosome with an LOH event, the copy number
594 across the whole chromosome is shown in the top part of the panel, with the total copy
595 number in black and the minor allele copy number in blue. Underneath are shown example
596 single nucleotide polymorphisms that support the LOH. In each case, reads from the crypt in
597 question are shown above, and reads from its matched normal are shown underneath. Thus in

598 the first image, the wild-type state (below) is heterozygous for a T (red) single nucleotide
599 polymorphism, whereas in the crypt in question (above), this polymorphism has now become
600 homozygous. Small deviations from a fully homozygous state are likely due to stromal
601 contamination. **(d)** reads supporting structural variants in normal colon are shown.

602

603 **Extended Data Figure 8. Gain of function driver mutations in normal colon.** Putative
604 driver missense mutations in oncogene hotspots. The number of substitutions catalogued in
605 COSMIC are shown on the y axis at each position along the gene, with the mutations
606 observed in our cohort highlighted.

607

608 **Extended Data Figure 9. Occurrence matrix of signatures and driver mutations in**
609 **crypts.** For all crypts that were whole genome sequenced to sufficient depth and for crypts
610 that underwent targeted sequencing and in which driver mutations were found, the signatures
611 and driver mutations are shown. Each vertical column represents a crypt. The individual to
612 which each crypt belongs is indicated by alternating colours in the top bar. The site to which
613 each crypt belongs is shown underneath. The matrix is coloured by contribution of each
614 signature to each crypt, normalised for each signature; thus the crypt with the largest
615 contribution of a given signature is coloured purple, and the crypt with the smallest
616 contribution is coloured white. Crypts in which the signatures could not be assessed, either
617 because they underwent targeted sequencing or the coverage was poor, are coloured grey.
618 Driver mutations, including heterozygous mutations in tumour suppressor genes, are
619 indicated by a black bar.

620

621 **Extended Data Figure 10. Stem cell dynamics of normal colon** (Supplementary Results 2).
622 **(a)** stem cell number and stem cell replacement rate in normal human colonic crypts, as
623 estimated by approximate Bayesian computation. Every point represents a simulation. Points
624 are coloured according to their similarity to the observed data: the most similar 0.1% are
625 coloured dark red, and so on, until the least similar simulations are blue. **(b)** approximate
626 Bayesian computation of the crypt fission rate in the human colon. The prior distribution of
627 the crypt fission rate used to simulate many biopsies of the colon is shown above. The unit
628 for the crypt fission rate is fissions per crypt per year. The posterior distribution of the crypt
629 fission rate parameter estimated by neural network regression on the simulations is shown
630 below. **(c-d)**, evidence of crypt fusion in human colon. In each figure, above is shown a
631 phylogeny that depicts the genetic relationships between selected crypts. Dotted blue lines
632 show low allele fraction mutations that are shared between crypts in a manner that is
633 incompatible with phylogeny that is dictated by the clonal mutations. Underneath each crypt
634 in the phylogeny is an image depicting its position in the section. Sections are labelled
635 according to their z-stacked order. Beneath this, the allele fraction of mutations on each
636 branch of the phylogeny in each crypt is shown. To the right hand side is the trinucleotide
637 context of the mutations that occurred on the branch.

638

639 **Extended Data Figure 11. Comparison of the mutational signatures and driver**
640 **landscape of normal crypts and colorectal adenocarcinomas.** **(a)** a comparison of the
641 burden of mutations due to every mutational signature found in either group. For each
642 signature, the (mutation burden+1) of every sample is shown on the y axis on a log scale.
643 Normal colon and cancer samples are ordered within their groups. Colorectal
644 adenocarcinoma signature attributions and burden are from Alexandrov et al.¹. 60 cancers are
645 compared with 472 normal crypts. **(b)** the proportion of driver mutations in each gene in
646 normal colon and colorectal cancer. The frequency of driver mutations in cancer was derived
647 using data from The Cancer Genome Atlas Network⁴³ (Supplementary Methods).

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697

SUPPLEMENTARY METHODS

Human tissues

We obtained healthy colonic biopsies from four cohorts (Supplementary Table 1). The first represents seven deceased organ donors ranging in age from 36 to 67, from whom colonic and small intestinal biopsies were taken at the time of organ donation (REC 15/EE/0152). The second represents individuals aged 60 to 72 who were having a colonoscopy following a positive faecal occult blood test as part of the Bowel Cancer Screening Programme (Ethical approval 08-H0308-13); we selected 16 who were not found to have either an adenoma or a carcinoma on colonoscopy, and 15 who were found to have a colorectal carcinoma (the normal biopsies that we use were distant from these lesions). The third cohort represents three paediatric patients who underwent routine colonoscopy to exclude inflammatory bowel disease and who were found to have a completely normal intestinal mucosa macroscopically and histologically (REC 12/EE/0482). The final cohort included one 78 year-old gentleman with oesophageal cancer who underwent a warm autopsy (REC 13/EE/0043). This gentleman had been treated with palliative chemotherapy of Epirubicin, Oxaliplatin and Capecitabine within the three months before the autopsy; given that monoclonal conversion within crypts is on the order of years, mutations due to these chemotherapies are likely to be private to a small proportion of stem cells per crypt and so are unlikely to be detected. All samples were obtained with informed consent and studies approved by East of England Research Ethics Committees.

Laser capture microdissection of colonic crypts

Fresh frozen biopsies were embedded in optimal cutting temperature (OCT) compound. 30 micrometre sections were fixed in methanol for five minutes, washed three times with phosphate-buffered saline, and stained with Gill's haematoxylin for 20 seconds. Crypts were isolated by laser capture microdissection, and collected in separate wells of a 96-well plate. They were lysed using the Arcturus PicoPure Kit (Applied Biosystems) according to the manufacturer's instructions. DNA library prep then proceeded without clean-up or quantification.

Library preparation

Two library preparation methods were used for laser capture microdissected (LCM) material: in initial experiments sonication was used to fragment DNA, and later, an enzymatic fragmentation method was implemented as it could make libraries from even lower input. Comparison of the two methods showed no difference in mutation calls once post-processing filters (described below) had been implemented. All samples in this study were processed using an Agilent Bravo Workstation (Option B; Agilent Technologies).

For sonication libraries, LCM lysate (20 µl) was mixed with 100 µl TE buffer (Ambion; 10 mM Tris-HCl, 1 mM EDTA) and DNA was fragmented using focused acoustics (Covaris LE220; Covaris, Inc.). Fragmented DNA was mixed with 80 µl Ampure XP beads (Beckman Coulter). Following a 5 min binding reaction and magnetic bead separation, genomic DNA was washed twice with 75% ethanol. Beads were resuspended in 20 µl nuclease-free water (Ambion) and processed immediately for DNA library construction. Each sample (20 µl) was mixed with 2.8 µl of NEBNext Ultra II End Prep Reaction Buffer, 1.25 µl of NEBNext Ultra

698 II End Prep Enzyme Mix (New England BioLabs) and incubated on a thermal cycler for 30
699 min at 20°C then 30 min at 65°C. Following DNA fragmentation and A-tailing, each sample
700 was incubated for 20 min at 20°C with a mixture of 30 µl ligation mix and 1 µl ligation
701 enhancer (New England BioLabs), 0.9 µl nuclease-free water (Ambion) and 0.1 µl duplexed
702 adapters (100 uM; 5'-ACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3', 5'-phos-
703 GATCGGAAGAGCGGTTCAGCAGGAATGCCGAG-3'). Adapter-ligated libraries were
704 purified using Ampure XP beads by addition of 65 µl Ampure XP solution (Beckman
705 Coulter) and 65 µl TE buffer (Ambion). Following elution and bead separation, DNA
706 libraries (21.5 µl) were amplified by PCR by addition of 25 µl KAPA HiFi HotStart
707 ReadyMix (KAPA Biosystems), 1 µl PE1.0 primer (100 µM; 5'-
708 AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGA
709 TC*T-3') and 2.5 µl iPCR-Tag (40 µM; 5'-
710 CAAGCAGAAGACGGCATAACGAGATXGAGATCGGTCTCGGCATTCCTGCTGAACC
711 GCTCTTCCGATC-3') where 'X' represents one of 96 unique 8-base indexes The sample
712 was then mixed and thermal cycled as follows: 98 °C for 5 min, then 12 cycles of 98 °C for
713 30 s, 65°C for 30 s, 72 °C for 1 min and finally 72 °C for 5 min. Amplified libraries were
714 purified using a 0.7:1 volumetric ratio of Ampure Beads (Beckman Coulter) to PCR product
715 and eluted into 25 µl of nuclease-free water (Ambion). DNA libraries were adjusted to 2.4
716 nM and sequenced on the HiSeq X platform (illumina) according to the manufacturer's
717 instructions with the exception that we used iPCRtagseq (5'-
718 AAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTC-3') to read the library index.

719

720 For enzymatic fragmentation, LCM lysate (20 ul) was mixed with 50 ul Ampure XP beads
721 (Beckman Coulter) and 50 µl TE buffer (Ambion; 10 mM Tris-HCl, 1 mM EDTA) at room
722 temperature. Following a 5 min binding reaction and magnetic bead separation, genomic
723 DNA was washed twice with 75% ethanol. Beads were resuspended in 26 µl TE buffer and
724 the bead/genomic DNA slurry was processed immediately for DNA library construction.
725 Each sample (26 µl) was mixed with 7 µl of 5X Ultra II FS buffer, 2 µl of Ultra II FS enzyme
726 (New England BioLabs) and incubated on a thermal cycler for 12 min at 37°C then 30 min at
727 65°C. Following DNA fragmentation and A-tailing, each sample was incubated for 20 min at
728 20°C with a mixture of 30 µl ligation mix and 1 µl ligation enhancer (New England
729 BioLabs), 0.9 µl nuclease-free water (Ambion) and 0.1 µl duplexed adapters (100 uM; 5'-
730 ACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3', 5'-phos-
731 GATCGGAAGAGCGGTTCAGCAGGAATGCCGAG-3'). Adapter-ligated libraries were
732 purified using Ampure XP beads by addition of 65 µl Ampure XP solution (Beckman
733 Coulter) and 65 µl TE buffer (Ambion). Following elution and bead separation, DNA
734 libraries (21.5 µl) were amplified by PCR by addition of 25 µl KAPA HiFi HotStart
735 ReadyMix (KAPA Biosystems), 1 µl PE1.0 primer (100 µM; 5'-
736 AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGA
737 TC*T-3') and 2.5 µl iPCR-Tag (40 µM; 5'-
738 CAAGCAGAAGACGGCATAACGAGATXGAGATCGGTCTCGGCATTCCTGCTGAACC
739 GCTCTTCCGATC-3') where 'X' represents one of 96 unique 8-base indexes The sample
740 was then mixed and thermal cycled as follows: 98 °C for 5 min, then 12 cycles of 98 °C for
741 30 s, 65°C for 30 s, 72 °C for 1 min and finally 72 °C for 5 min. Amplified libraries were
742 purified using a 0.7:1 volumetric ratio of Ampure Beads (Beckman Coulter) to PCR product
743 and eluted into 25 µl of nuclease-free water (Ambion). DNA libraries were adjusted to 2.4
744 nM and sequenced on the HiSeq X platform (Illumina) according to the manufacturer's
745 instructions with the exception that we used iPCRtagseq (5'-
746 AAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTC-3') to read the library index.

747

748 **Whole genome sequencing**

749 We generated paired end sequencing reads (150bp) using Illumina XTEN® machines
750 resulting in ~15x coverage per sample. In 94% of the whole genome crypts included for
751 statistical analysis, over 90% of the callable genome was covered by more than 10 reads.
752 Sequences were aligned to the human reference genome (NCBI build37) using BWA-MEM.

753

754 **Targeted sequencing**

755 A 2.3 MB capture panel was designed in-house to pull down genes that are known or
756 suspected to play a role in neoplasia. We performed custom RNA bait design following the
757 manufacturer's guidelines (SureSelect, Agilent). Samples were multiplexed on flow cells and
758 subjected to paired end sequencing (75-bp reads) using Illumina HiSeq2000 machines. One
759 96-well plate of samples was sequenced on each lane, but as tissue recovery was variable, a
760 range of coverage was achieved. Sequences were aligned to the human reference genome
761 (NCBI build37) using BWA-align.

762

763 **Data Availability**

764 Whole genome and targeted sequencing data are deposited in the European Genome
765 Phenome Archive (EGA) with EGA accession EGAD00001004192 and EGAD00001004193.
766 Images of microdissections and the physical distances between crypts are available on
767 Mendeley Data by searching for the title of this article. All other data is available from the
768 authors on request.

769

770 **Code Availability**

771 Code for statistical analyses is provided as part of the supplement. Custom R scripts and their
772 input data for signature analysis are available on GitHub at https://github.com/HLee-Six/colon_microbiopsies. All other code is available from the authors on request.

773

774 **Calling substitutions**

775 Substitution calling was broken down into three steps: mutation discovery; filtering to
776 produce a list of clean sites; and genotyping, where the presence or absence of every
777 mutation in every sample is evaluated.

778

779 First, mutations were initially discovered using the Cancer Variants through Expectation
780 Maximisation (CaVEMan) algorithm⁴⁶. CaVEMan uses a naïve Bayesian classifier to derive
781 the probability of all possible genotypes at each nucleotide. CaVEMan copy number options
782 were set to major copy number 5 and minor copy number 2 for normal clones, as in our
783 experience this maximises sensitivity. The algorithm was run using an unmatched normal in
784 order to be able to derive phylogenies: had another sample from the same individual been
785 treated as a matched normal, early embryonic mutations would have been treated as germline
786 and discarded, resulting in incorrect trees.

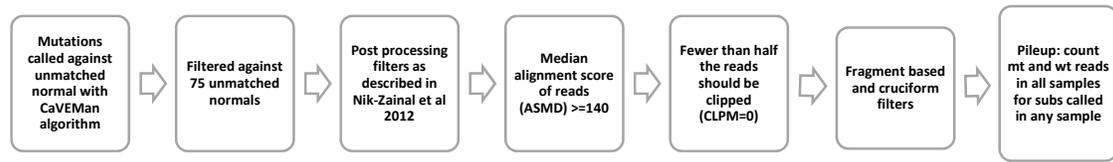
787

788 Second, a number of post-processing filters were applied (Extended Data Figure 2). These
789 included filtering against a panel of 75 unmatched normal samples to remove common single
790 nucleotide polymorphisms, post-processing as described previously³² and two filters (only
791 applied to whole genome sequencing data) designed to remove mapping artefacts associated
792 with BWA-MEM: the median alignment score of reads supporting a mutation should be
793 greater than or equal to 140, and fewer than half of these reads should be clipped. The library
794 preparation protocol for microbiopsies produced shorter library insert sizes than standard
795 methods. Reads could therefore overlap, resulting in double counting of mutant reads.
796 Fragment-based statistics were generated to prevent the calling of variant supported by a low
797

798 number of fragments. Variants were annotated by ANNOVAR⁴⁷ and fragment-based
799 statistics (fragment coverage, number of fragments supporting the variant, fragment-based
800 allele fraction) were calculated for each variant after the exclusion of marked PCR duplicates.
801 In the rare event of discordance in the called base at the variant position between overlapping
802 paired-end reads, the base with the highest quality score was selected. Fragment-based
803 statistics were calculated separately for high quality fragments (alignment score ≥ 40 and
804 base scores ≥ 30). Variants supported by at least three high quality fragments were retained
805 and used for the next stage of variant filtering. Inspection of variants specific to LCM
806 experiments revealed that the vast majority were present within inverted repeats capable of
807 forming hairpin structures, that they were supported by reads with very similar alignment
808 start position (and so not marked as PCR duplicates), and were primarily located close to the
809 alignment start within the supporting reads. Commonly these variants coincided with other
810 proximal variants (1-30 bp), but filtering based on variant proximity would also remove
811 actual kataegis events. *In silico* modelling of the potential hairpin showed that the variants
812 were aligning to each other in the stem of the structure, but could not form a base pair, while
813 all other bases could. The artefacts are likely the consequence of erroneous processing of
814 cruciform DNA (existing either prior to DNA isolation or formed during library preparation)
815 by the enzymatic digestion protocol applied. We have considered modelling the hairpin
816 structures to filter these variants, but given the fact that read clustering (i.e., similar alignment
817 position) serves as a hallmark for these artefacts, we opted to use the proximity of the variant
818 to the alignment start, and the standard deviation (SD) and median absolute deviation (MAD)
819 of the variant position within the supporting reads, as features for filtering. These statistics
820 were calculated separately for positive and negative strand aligned reads. In case the variant
821 was supported by a low number of reads (i.e., 0-1 reads) for one of the strands, the filtering
822 was based only on the statistics generated for the other strand. Per variant, if one of the
823 strands had too few reads supporting, it was required for the other strand that either: (I) there
824 should be $\leq 90\%$ supporting reads to report the variant within the first 15% of the read
825 starting from the alignment start, or (II) the statistics $MAD > 0$ and $SD > 4$. Per variant, if
826 both strands were supported by sufficient reads it was required for both strands separately
827 that either: (I) there should be $\leq 90\%$ supporting reads to report the variant within the first
828 15% of the read, (II) the statistics $MAD > 2$ and a $SD > 2$, or (III) that the other strand should
829 have the statistics $MAD > 1$ and $SD > 10$ (i.e., the variant is retained if the other strand
830 demonstrates strong measures of variance). In our experience, the proposed strategy vastly
831 reduces the number of artefactual variants while retaining all other variants, as assessed by
832 running the last filtering step on WGS data from non-LCM experiments.

833
834 Third, mutations were genotyped in every sample. A pileup of all the samples from a given
835 individual was constructed, counting the number of mutant and wild type reads in every
836 sample over every site that had been called in any sample from that person. Only reads with a
837 mapping quality of 30 or above and bases with a base quality of 30 or above were counted.
838 After applying these filters, mutations were genotyped based on the number of mutant and
839 wild type reads at each locus. Mutations were called based on a variant allele fraction (VAF)
840 > 0.2 , a depth > 7 , and at least 4 mutant reads. If the depth over a locus was less than seven in
841 a given sample, or if there was more than one mutant read but the other criteria were not met,
842 the genotype was set to NA for tree construction purposes. Loci that were set to NA in more
843 than one third of the samples were removed for construction of the phylogeny. Positions were
844 called as germline if they were either called as present or NA in all of the samples from a
845 given individual.
846

847 1.2% of all mutations were present in the coding regions of the genome. All mutations in
848 coding regions are provided (Supplementary Table 3).
849



850
851

852 **Calling short insertions and deletions (indels)**

853 As for substitutions, calling of indels was broken down into mutation discovery, filtering, and
854 genotyping. Mutations were called with the Pindel algorithm⁴⁸ using an unmatched normal.
855 Post processing filters were applied as in Nik-Zainal et al.³², and the number of mutant and
856 wild-type reads was tabulated as above. The same dataset-specific filters were applied as for
857 substitutions. Indels were then genotyped based on a VAF>0.2, a depth of at least 10, and
858 support of at least 5 mutant reads.

859

860 **Calling structural variants**

861 Genomic rearrangements were called using the BRASS algorithm
862 (<https://github.com/cancerit/BRASS>). Abnormally paired read pairs from WGS were grouped
863 and filtered by read remapping. Read pair clusters with $\geq 50\%$ of the reads mapping to
864 microbial sequences were removed, as were rearrangements where the breakpoint could not
865 be reassembled. Candidate breakpoints were matched to copy number breakpoints defined by
866 ASCAT (see below) within 10kb. Only structural variants where the two breakpoints were
867 more than 1000 base pairs apart were considered. Structural variants were called against a
868 matched normal skin or blood sample when available and against another crypt from the
869 same individual with good coverage when not.

870

871 **Calling copy number**

872 Copy number changes were called using the Allele-Specific Copy number Analysis of
873 Tumours (ASCAT) algorithm⁴⁹. The same matched normal sample was used as for calling
874 structural variants. For additional validation of copy number changes in normal colon, the
875 QDNAseq algorithm⁵⁰ was run. ASCAT uses both the read depth and ratios of heterozygous
876 single nucleotide polymorphisms to determine an allele-specific copy number, while the
877 QDNAseq relies solely on variations in sequencing coverage. To call amplifications and
878 deletions in the colonic microbiopsy cohort, only those that were both called by ASCAT and
879 showed a clear departure from the background \log_2 ratio by QDNAseq were retained. To call
880 copy neutral loss of heterozygosity in this cohort, all such events called by ASCAT were
881 checked visually on Jbrowse⁵¹ to verify an imbalance of parental snps. Only crypts with
882 >10X coverage, for which copy number changes could be reliably detected, were used.

883

884 **Detection of driver variants and positive selection**

885 Driver mutations were detected both through an unbiased dNdS method and through manual
886 annotation. For these analyses, the CaVEMan and Pindel calls were used without post-
887 processing filters (such as requiring a VAF cutoff of >0.2) in order to maximise our
888 sensitivity. All putative driver variants were visually inspected using Jbrowse⁵¹, and so we
889 could afford a higher false positive rate in the mutation discovery phase.

890

891 dNdScv⁵² was used to conduct three tests: first, using only the whole genome sequencing
892 data, an analysis of selection over all genes; second, using combined whole genome and
893 targeted sequencing data, over all the genes covered by the bait-set; and finally, using again
894 this combined dataset, over 90 selected cancer genes (Supplementary Table 4, Supplementary
895 Results 1).

896
897 Manual annotation of driver variants based on prior knowledge complemented this. A list of
898 90 colorectal cancer genes (appendix) curated from the literature that were also covered by
899 the bait-set were intersected with the list of substitutions and indels from combined whole
900 genome and targeted sequencing. Mutations were annotated as putative drivers if they were
901 either missense mutations that fell in an oncogene hotspot (based on visualisation of the
902 distribution of mutations in the gene on COSMIC⁵³), or if they were truncating mutations that
903 fell in a tumour suppressor gene.

904
905 Structural variants that might act as drivers were assessed by intersection of genes involved
906 in each structural variant with the twelve genes involved in gene fusions that have been
907 reported in colorectal cancer in COSMIC (*VTI1A*, *TCF7L2*, *TPM3*, *NTRK1*, *PTPRK*, *RSPO3*,
908 *ETV6*, *NTRK3*, *EIF3E*, *RSPO2*, *C2orf44*, and *ALK*). No fusion genes were found. None of the
909 genes involved in structural variants in our data overlapped with the list of 90 cancer genes
910 used for assessing substitutions and indels, and nor were there any genes that were affected
911 by more than one structural variant. No high-level copy number amplifications were observed
912 and there were no homozygous deletions.

913
914 Please note that the driver frequency is low, such that we cannot estimate a per-gene driver
915 frequency. All we can do to derive a meaningful estimate is to pool our driver mutations. In
916 addition, coverage may fluctuate even within a gene. Some portions of a gene may be well
917 covered in 1,000 crypts, and others in 2,000 crypts. The approach that we took was to
918 calculate, for the average exonic base pair in our 90 cancer genes, the number of crypts in
919 which that base pair was covered by ≥ 8 reads (for substitutions) and by ≥ 10 reads (for
920 indels). 64% of all bases in the targeted panel across all crypts are covered by ≥ 8 reads,
921 which equates to a number of callable bases equivalent to having sequenced $\sim 1,400$ crypts
922 with perfect coverage over every base in every crypt. This average number of crypts in which
923 all base pairs achieve good coverage becomes the denominator for calculating the driver
924 mutation frequency (with the number of drivers observed in the dataset as the numerator). A
925 similar approach can be taken with indels. Our estimate of 1% uses a global correction, on
926 the assumption that under-representation and over-representation will even itself out when
927 estimating the total frequency of driver mutations in the whole dataset.

928

929 **Estimation of frequency of driver mutations in cancer**

930 Publicly-available colorectal cancer mutation calls were obtained from The Cancer Atlas
931 Network³⁷. Driver mutations were annotated manually in the same way as in our dataset: only
932 mutations that fell in the 90 genes that we had selected were considered, and they were
933 annotated as putative drivers if they were either missense mutations that fell in an oncogene
934 hotspot (based on visualisation of the distribution of mutations in the gene on COSMIC⁵³), or
935 if they were truncating mutations that fell in a tumour suppressor gene.

936

937 **Construction of phylogenies**

938 Phylogenies are used in this analysis for timing mutations. The most informative branches in
939 this case are the long branches shared by a small number of crypts, which are very robust to
940 all tree construction methods. Trees were built using maximum parsimony using substitutions

941 called as described above. For every individual, the input matrix of mutation calls was
942 bootstrapped 100 times. Phylogenies were constructed for each replicate using the Wagner
943 method of the Mix programme from the Phylip suite of tools⁵⁴. The consensus phylogeny was
944 constructed from 100 bootstrap runs using the extended majority rule method for the
945 Consense programme from the Phylip suite of tools⁵⁴.

946 Across all phylogenies, a mean of 10% and a median of 1.5% of mutations per tree
947 did not fit the trees perfectly. Phylogenies with more crypts had more mutations that fitted
948 imperfectly. Consider a mutation that is really present in 50 crypts. Even with 15X coverage
949 over the site in every sample, and with every crypt completely clonal, if we simulate
950 resampling of mutant reads from the binomial distribution (with size of 15 and probability of
951 0.5), 17% of the time the mutation will have fewer than the 3 reads required to call it in at
952 least one sample. Variation in sequencing depth, clonality, and sequencing errors would
953 further decrease the probability of calling the mutation perfectly in every sample. Nodes
954 across all our phylogenies had mean bootstrapping values of 0.77 and median bootstrapping
955 values of 0.99. Branches at the very top of the phylogenies, likely representing embryonic
956 cell divisions, are supported by only a few mutations and have lower support because in a
957 given bootstrap sample the couple of mutations that support this node may be omitted.
958 Longer shared branches almost always have bootstrapping values of 1. These longer shared
959 branches are those that are most important to our analyses, because they are the most
960 informative when timing mutational signatures relative to one another and because they
961 represent postnatal crypt fission events. In order to increase further our confidence in our
962 phylogenies, we validated them by reconstructing them with indels. To do this, the same
963 procedure as for substitutions was followed for indel matrices. As there were fewer indels
964 than substitutions, nodes in indel phylogenies were generally reconstructed with lower
965 confidence than in substitution phylogenies, but they broadly agree. 85% of nodes
966 reconstructed with $\geq 90\%$ confidence in the indel tree were present with exactly the same set
967 of descendants in the substitution trees. Any errors in the phylogenies should be relatively
968 minor and not affect our downstream analyses.

969
970 The phylogeny inference programme used provided the topology of the tree but not the
971 assignment of mutations. Mutations from the input matrix of genotypes therefore have to be
972 re-assigned to branches. In order to assign a set of mutation calls with no false negative and
973 no false positives to a tree, each branch of the tree was considered in turn. If a mutation was
974 called in all the descendants of a given branch, and in no samples that were not descendants
975 of the branch, mutations were assigned to that branch.

976
977 Some colonic microbiopsies suffered from low coverage and stromal contamination. For this
978 reason, we did not expect mutations to fit the tree perfectly, as a mutation that was truly
979 present in a colony might be missed if too few supporting reads are found. Mutations were
980 only assigned to the tree in order to determine the mutational processes active at a particular
981 time. We reasoned that it was preferable to assign only mutations that fit the tree perfectly
982 and adjust the branch lengths based on the power to call mutations at a given branch, rather
983 than attempting to assign mutations that fit the tree imperfectly. Using the clonality and
984 coverage of all descendants of a branch, the proportion of true substitutions or indels on the
985 branch that would be first discovered (whether by CaVEMan or Pindel) and then genotyped
986 as present according to the criteria described above was calculated. The observed branch
987 length was then adjusted by dividing by this proportion. Adjustment proportions can be found
988 in Supplementary Table 7. This was done for both substitutions and indels, but not for
989 structural variants and for larger copy number changes due to a lack of data: most branches

990 have no large variants and so could not be extended appropriately. Rearrangements and copy
991 number changes were assigned to phylogenies manually.

992

993

994 **Extraction of mutational signatures**

995 Mutational signatures were extracted using the mutations assigned to every branch of a
996 phylogeny as a ‘sample’. This allows better discrimination of mutational processes that may
997 occur at different times within the same cell. Mutations were categorised following the
998 method used by the Mutational Signatures working group of the Pan Cancer Analysis of
999 Whole Genomes (PCAWG)¹. Single base substitutions were categorised into 96 classes
1000 according to the identity of the pyrimidine mutated base pair, and the base 5’ and 3’ to it.
1001 Doublet base substitutions were categorised into 78 classes according to the identity of the
1002 reference and alternative bases. Indels were classified according to whether they were an
1003 insertion or a deletion, the identity of the inserted/deleted base, the length of the
1004 mononucleotide tract in which they occurred, or the degree of homology with the
1005 surrounding sequence into 83 classes (Fig. 1a).

1006

1007 Signatures were extracted using a hierarchical Dirichlet Process^{55,56}. Code and the input
1008 mutations are provided at https://github.com/HLee-Six/colon_microbiopsies. First, the
1009 algorithm was conditioned on the set of mutational signatures that have found to be operative
1010 in colorectal cancers in PCAWG¹: SBS1, SBS2, SBS3, SBS5, SBS13, SBS16, SBS17a,
1011 SBS17b, SBS18, SBS25 (included although it is not found in colorectal cancer because the
1012 similarity with the mutational profile with crypts from one individual had been previously
1013 noted), SBS28, SBS30, SBS37, SBS40, SBS41, SBS43, SBS45, SBS49, DBS, DBS3, DBS4,
1014 DBS6, DBS7, DBS8, DBS9, DBS10, DBS11, ID1, ID2, ID3, ID4, ID5, ID6, ID7, ID8, ID10,
1015 and ID14. This allows simultaneous discovery of new signatures and matching to known
1016 ones. Nine single base substitution (SBS), two doublet base substitution (DBS), and five
1017 indel (ID) signatures were discovered (Extended Data Fig. 2). Despite pre-conditioning,
1018 signatures that were perfectly correlated in all samples were still amalgamated. This
1019 occurred, for example, with signatures 1, 5, and 18. Therefore, expectation maximisation was
1020 used to deconvolute all HDP signatures into known PCAWG signatures. If a signature
1021 reconstituted from the components that expectation maximisation extracted (only including
1022 PCAWG signatures that accounted for at least 10% of mutations in each sample to avoid
1023 over-fitting) had a cosine similarity to the HDP signature of more than 0.95, the signature
1024 was presented as its expectation maximisation deconvolution. Three HDP signatures met
1025 these criteria: the HDP SBS1 signature was deconvoluted into a mixture of PCAWG SBS1,
1026 PCAWG SBS5, and PCAWG SBS18; the HDP DBSA was deconvoluted in PCAWG DBS2,
1027 PCAWG DBS4, PCAWG DBS6, PCAWG DBS9, and PCAWG DBS11; and the HDP IDC
1028 was deconvoluted into PCAWG ID1, PCAWG ID2, and PCAWG ID5 (Extended Data Fig.
1029 3). To test the robustness of this signature analysis, other signature extraction methods were
1030 used: HDP with no pre-conditioning, the non-negative matrix factorisation (NNMF) method
1031 used by Blokzijl and colleagues⁴, and a version of the NNMF algorithm used by Alexandrov
1032 and colleagues¹. These all produced comparable results (Extended Data fig. 4).

1033

1034 **Timing SBSA and SBSB throughout life**

1035 Five patients had informative clades with branchpoints that allowed us to time SBSA.
1036 Plotting the cumulative amount of SBSA vs SBS1 at each node in these clades (Extended
1037 Data Fig. 6aq), we observed that for each the rate of accumulation of SBSA relative to SBS1

1038 was high in early branchpoints, and then slowed down almost to zero on all branches but for
1039 one (a branch of patient ao, where it continues to be acquired, albeit at a slow rate).

1040

1041 We can take the inflexion point on the graph of cumulative SBSA vs SBS1 to be the upper
1042 limit of the point in time when SBSA slowed down. This provides an upper bound because:

1043

1044 1. When we observe the presence of a signature on a branch, we know that the causative
1045 process must have been active at some point during the lifetime of the branch, but we
1046 cannot say when on the branch it occurs; it might have ended long before the branch
1047 did.

1048 2. If the time to the most recent common ancestor of the crypt is longer than 0, the age at
1049 which this stopped would be earlier.

1050 3. If the SBS1 mutation rate is increased in early life, as it may be during the rapid
1051 growth of the embryo, the age at which the inflexion point occurs would have been
1052 younger.

1053

1054 Using these five informative clades, and assuming a clock-like but personalised rate of SBS1
1055 accumulation (i.e. each patient can accumulate SBS1 at their own constant rate), we found
1056 that the upper bound of the age at which SBSA slowed was: 9.7 years (patient h); 7.1 years
1057 (patient z); 2.4 years (patient am); 20.1 years (patient aa); and 9.6 years (patient ao). There
1058 are no branches that begin after 10 years of age with a high ratio of SBSA to SBS1.

1059

1060 The most informative branchpoint is the earliest inflexion point; the estimate of 2.4 years
1061 from patient aa, is therefore, perhaps our best estimate. Nonetheless, we did not want to base
1062 our statement on a single patient, and so 10 years was given in the text as four patients had
1063 branches that ended before 10 years of age.

1064

1065 A similar argument can be made for SBSB (Extended Data Fig. XXX). For SBSB, however,
1066 only two clades were informative. The estimated upper bounds of age for SBSB activity were
1067 2.4 years old (this was the same inflexion point as for SBSA in patient aa), and 6.4 years old
1068 (patient ai). For patient ai, a reasonable amount of SBSB is still acquired after this
1069 branchpoint. If the ratio of accumulation of SBSB vs SBS1 continued at the same rate as
1070 before this branchpoint, the number of SBSB mutations seen in the terminal branches would
1071 have been observed by age 8.2 years old.

1072

1073 **Telomere length analysis**

1074 Telomere length was estimated from whole genome sequencing data using Telomerecat.
1075 Telomerecat is a python based software package for telomere length estimation from short-
1076 read whole genome sequence data⁵⁷. It functions by classifying paired-end reads as either
1077 fully or partially telomeric based on the canonical hexamer TTAGGG, and uses that ratio to
1078 estimate an average telomere length. Notably, Telomerecat measures telomeres and also
1079 accounts for interstitial telomeric repeats. It is ploidy and species agnostic (assuming that the
1080 telomere hexamer is the canonical mammalian signature of TTAGGG_n). Telomerecat has
1081 four main stages; 1) Identification of all telomeric or partially telomeric read pairs and
1082 creation of a subsetted bam file containing only these reads 2) Classification of telomeric
1083 read-pairs into intratelomeric, boundary or junction-spanning or intrachromosomal 3) Error
1084 correction of boundary or junction-spanning read pairs 4) Estimation of telomere length
1085 based on the ratio of intratelomeric and boundary / junction spanning read pairs.

1086

1087 Telomerecat has been validated on whole genome DNA sequencing files from both tumour
1088 and normal samples⁵⁷. Its results show concordance with an established method of telomere
1089 length measurement; the mean Telomere Restriction Fragment (mTRF) technique.
1090 Alternative packages are available notably: Computel⁵⁸, Telseq⁵⁹ and TelomereHunter⁶⁰.
1091 They all have respective strengths and have been benchmarked through their methods
1092 publications. We have opted for Telomerecat as it provides a base pair resolution estimate of
1093 telomere length whilst providing correction for variations in sequencing depth in a ploidy
1094 agnostic manner.

1095

1096 We have run Telomerecat on 445 crypt bam files with good coverage and clonality to
1097 generate telomere length estimates. Telomerecat was threaded across 10 cores and 100
1098 simulation cycles were requested per run. Values displayed are the median telomere length
1099 across all chromosomes in that samples measured in base pairs.

1100

1101 **Statistical analyses**

1102 All statistical analyses were performed in R (Supplementary Results 1, Supplementary
1103 Results 2). Code can be found at https://github.com/HLee-Six/colon_microbiopsies.

1104

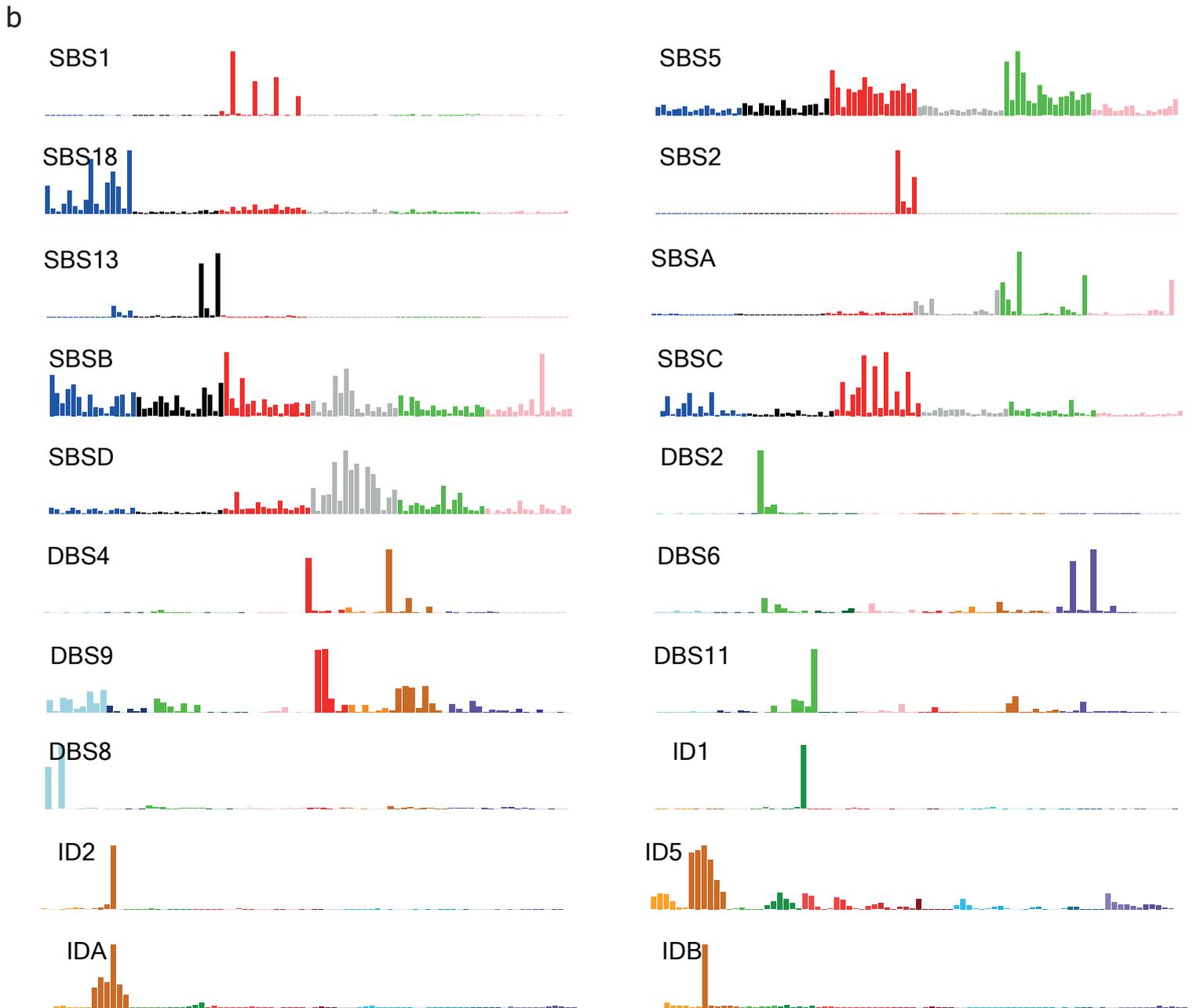
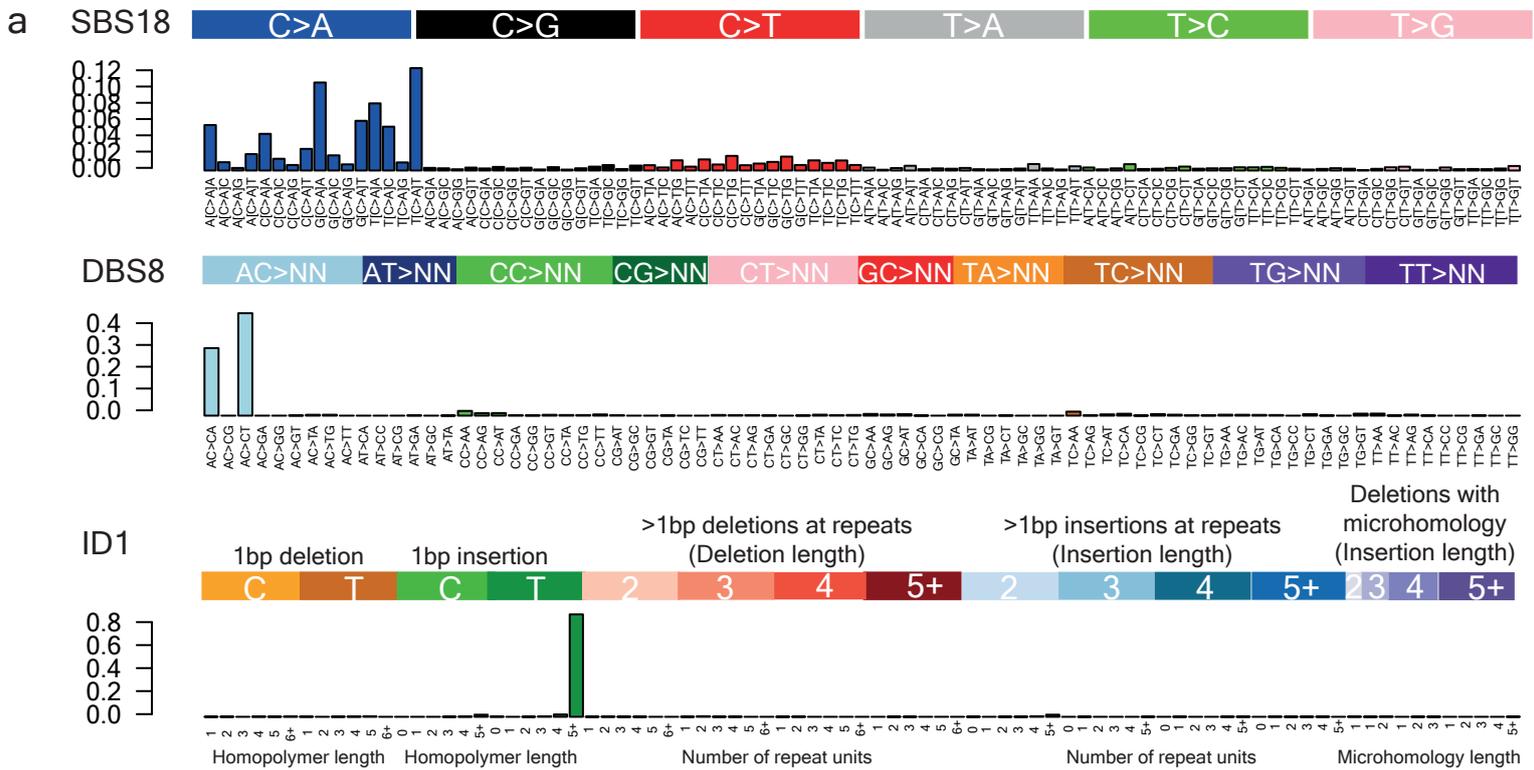
1105

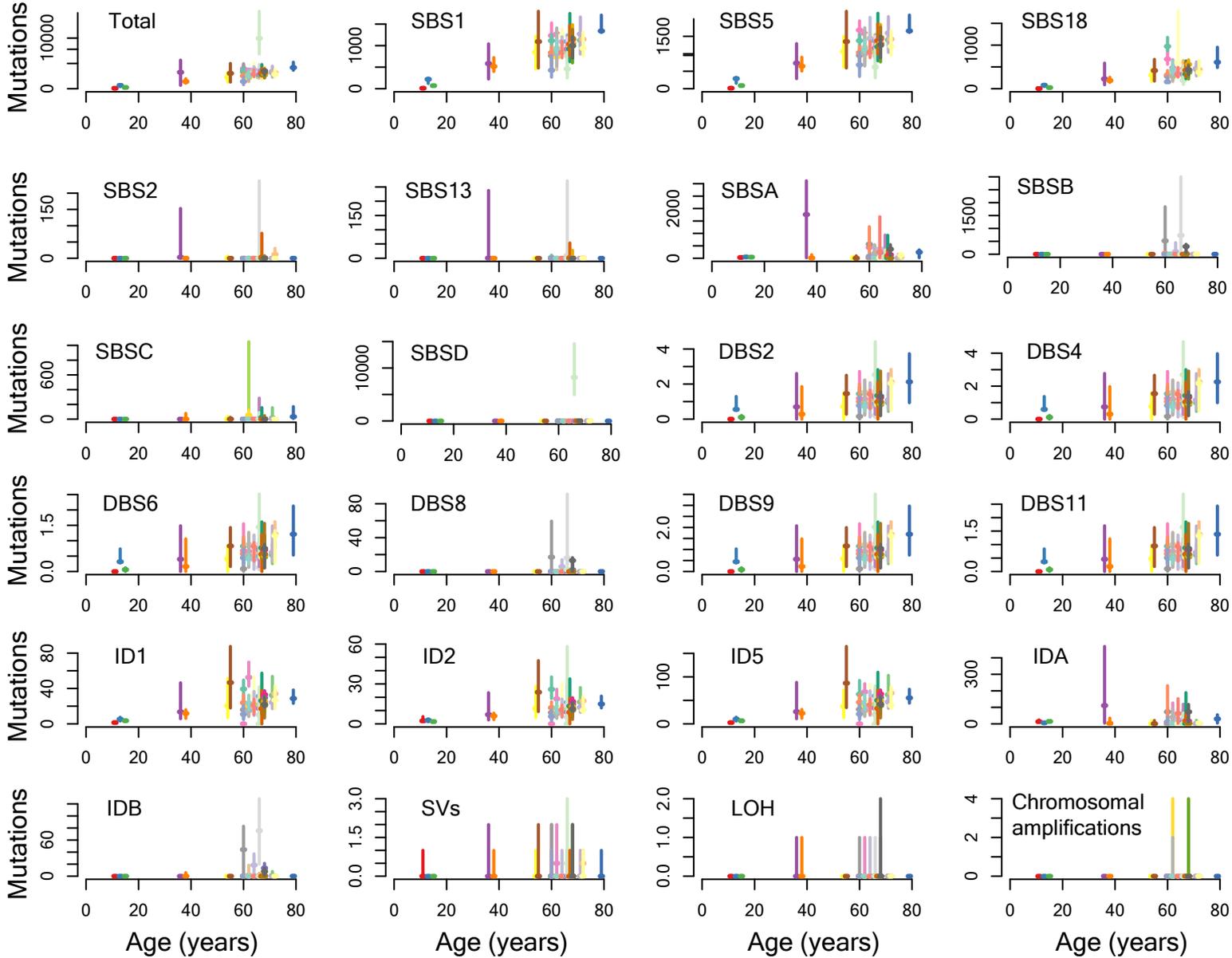
1106 **REFERENCES FOR SUPPLEMENTARY METHODS**

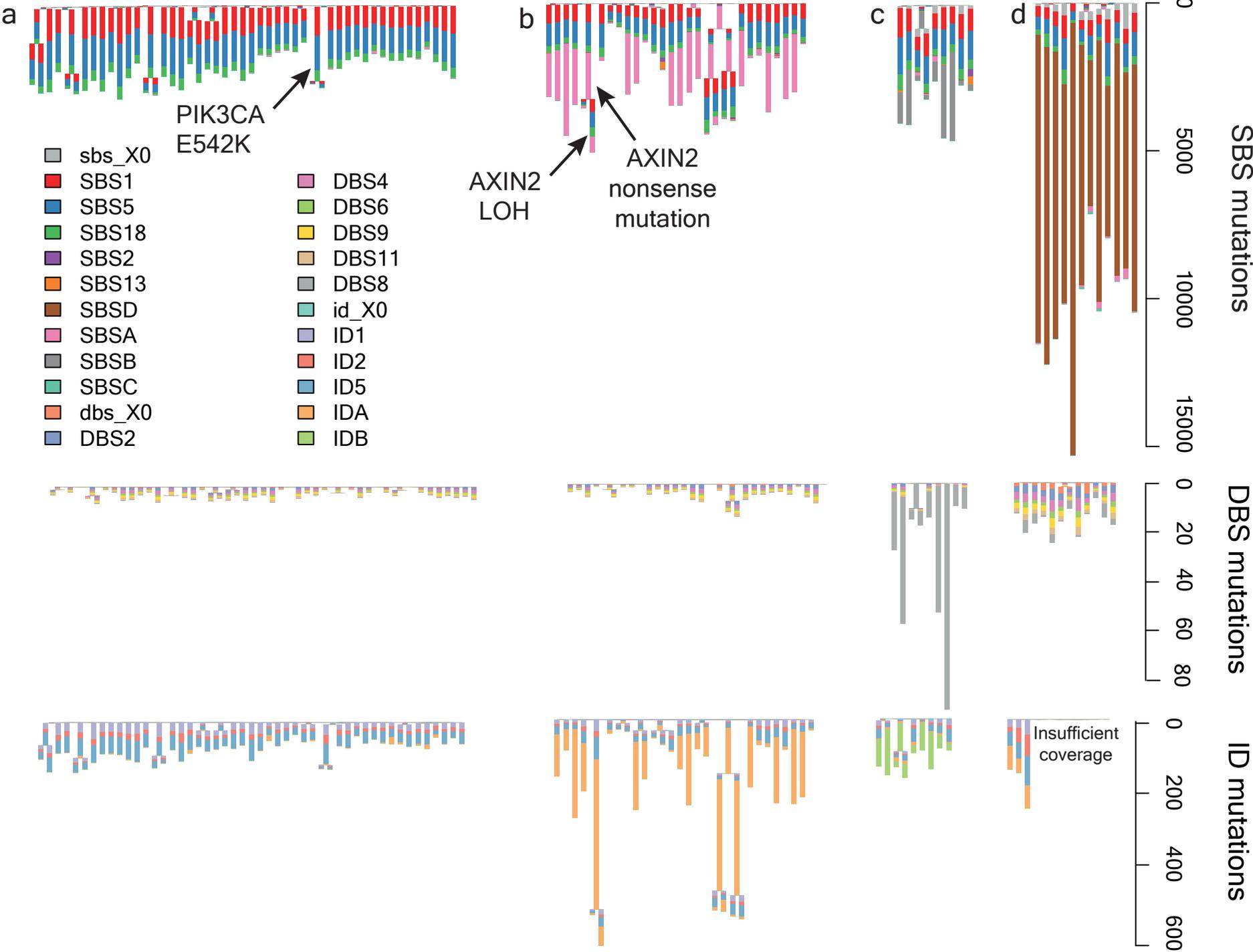
1107

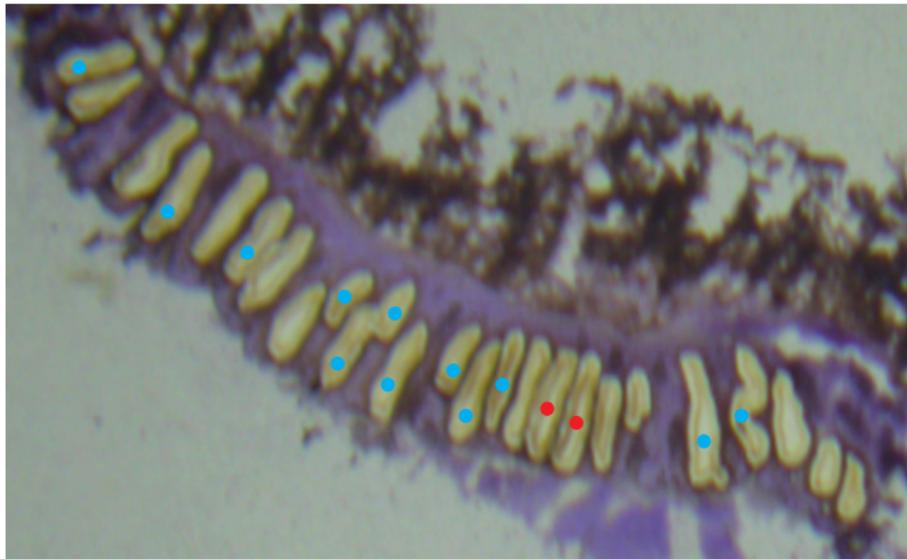
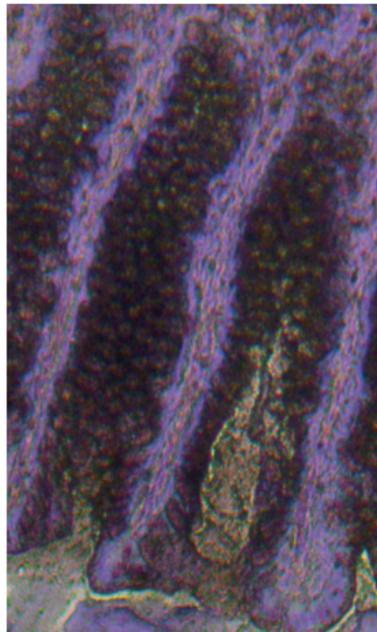
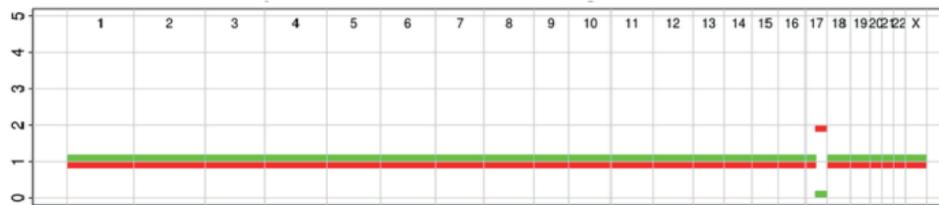
- 1108 46. Jones, D. et al. cgpCaVEManWrapper: Simple execution of CaVEMan in order to
1109 detect somatic single nucleotide variants in NGS data. *Curr. Protoc. Bioinformatics*
1110 **56**, 15.10.1–15.10.18 (2016).
- 1111 47. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic
1112 variants from high-throughput sequencing data. *Nucleic acids research* **38**, e164,
1113 (2010).
- 1114 48. Raine, K. M. et al. cgpPindel: identifying somatically acquired insertion and deletion
1115 events from paired end sequencing. *Curr. Protoc. Bioinformatics* **52**, 15.7.1–15.7.12
1116 (2015).
- 1117 49. Van Loo, P. et al. Allele-specific copy number analysis of tumors. *Proc. Natl Acad.*
1118 *Sci. USA* **107**, 16910–16915 (2010).
- 1119 50. Scheinin, I. et al. DNA copy number analysis of fresh and formalin-fixed specimens
1120 by shallow whole-genome sequencing with identification and exclusion of
1121 problematic regions in the genome assembly. *Genome Research*, **24**, 2022–2032
1122 (2014).
- 1123 51. Buels R *et al.* JBrowse: a dynamic web platform for genome visualization and
1124 analysis. *Genome Biology* doi: 10.1186/s13059-016-0924-1 (2016).
- 1125 52. Martincorena, I. et al. Universal patterns of selection in cancer and somatic tissues.
1126 *Cell* **171**, 1029–1041 (2017).
- 1127 53. Forbes S.A. et al. COSMIC: somatic cancer genetic sat high-resolution. *Nucleic Acids*
1128 *Res.* **45**, D777-D783 (2017).
- 1129 54. Felsenstein, J. PHYLIP — Phylogeny Inference Package (Version 3.2). *Cladistics* **5**,
1130 164–166 (1989).
- 1131 55. Roberts, N. *Patterns of somatic genome rearrangement in human cancer*. PhD thesis,
1132 Univ Cambridge, UK (Wellcome Trust Sanger Institute, 2018).
- 1133 56. Nicola Roberts, R pkg for Hierarchical Dirichlet Process,
1134 <https://github.com/nicolaroberts/hdp> (Accessed August 2018).

- 1135 57. Farmery, J. H. R., Smith, M. L., Bioresource, N., Diseases, R. & Lynch, A. G.
1136 Telomerecat : A ploidy-agnostic method for estimating telomere length from whole
1137 genome sequencing data. 1–17 (2018).
- 1138 58. Nersisyan, L. & Arakelyan, A. Computel: Computation of mean telomere length from
1139 whole-genome next-generation sequencing data. *PLoS One* **10**, 1–14 (2015).
- 1140 59. Ding, Z. *et al.* Estimating telomere length from whole genome sequence data. **42**, 7–
1141 10 (2018).
- 1142 60. Feuerbach, L., Sieverling, L., Deeg, K. I., Ginsbach, P. & Hutter, B. TelomereHunter :
1143 telomere content estimation and characterization from whole genome sequencing
1144 data. (2016).







a**b****c****d**