# The landscape of somatic mutation in normal colorectal epithelial cells
— **Source link** ↗

Henry Lee-Six, Sigurgeir Olafsson, Peter D. Ellis, Robert J. Osborne ...+22 more authors
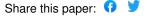
**Institutions:** Wellcome Trust Sanger Institute, Erasmus University Medical Center, University of Cambridge, Cambridge University Hospitals NHS Foundation Trust ...+1 more institutions

Related papers:

- Somatic mutant clones colonize the human esophagus with age

- High burden and pervasive positive selection of somatic mutations in normal human skin

- The Repertoire of Mutational Signatures in Human Cancer

- The mutational landscape of normal human endometrial epithelium

- Age-related remodelling of oesophageal epithelia by mutated cancer drivers

# The landscape of somatic mutation in normal colorectal epithelial cells

Henry Lee-Six[1], Peter Ellis[1], Robert J. Osborne[1], Mathijs A. Sanders[1,2], Luiza Moore[1], Nikitas Georgakopoulos[3], Franco Torrente[4], Ayesha Noorani[5], Martin Goddard[6], Philip Robinson[1], Tim H. H. Coorens[1], Laura O'Neill[1], Christopher Alder[1], Jingwei Wang[1], Rebecca C. Fitzgerald[5], Matthias Zilbauer[4], Nicholas Coleman[7], Kourosh Saeb-Parsy[3], Inigo Martincorena[1], Peter J. Campbell[1], Michael R. Stratton[1]*

1. Wellcome Trust 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, Cambridge University Hospital Trust, Addenbrookes, 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. 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 because of the difficulty of detecting somatic mutations in normal cells. Each colorectal crypt is a small clone of cells derived from a single recently-existing stem cell. 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 some phases of the cell lineage from zygote to adult cell. 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.

## 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, in substantial part, on identifying somatic mutations in morphologically abnormal neoplastic proliferations representing intermediate stages between normal and cancer cells. Classical studies of driver mutations in colorectal adenomas and carcinomas have been particularly influential in shaping our perspective in this regard[3].

53    As for most cancer types, however, the earliest stages of progression to colorectal cancer
54    remain considerably less well understood. The driver mutation that first sets a colorectal
55    epithelial cell on the path to cancer is likely caused by mutational processes operative in
56    normal cells, of which there is limited understanding. The nature and numbers of the earliest
57    neoplastic clones with driver mutations, which conceivably are morphologically
58    indistinguishable from normal cells, are similarly unclear. In large part, these deficiencies are
59    due to the technical challenge of identifying somatic mutations in normal tissues, which are
60    composed of myriad microscopic cell clones. Several approaches have been adopted to
61    address this, including sequencing of *in vitro* expanded cell populations derived from single
62    cells[4-9], sequencing normal tissue microbiopsies incorporating small numbers of clones[10,11],
63    sequencing single normal cells[12-14], highly error corrected sequencing[15], and non-sequencing
64    based approaches[27,44].
65
66    These approaches have provided insights into early stages of cancer development. Signatures
67    of common somatic mutational processes have been found in normal cells of the small and
68    large intestine, liver, blood, skin, and nervous system but thus far studies have not been of
69    sufficient scale to characterise variation in their activity or detect less frequent processes[4-10,
70    12-15]. Remarkably high proportions of normal skin epithelial cells have been shown to be
71    members of clones already carrying driver mutations[10], and large mutant clones have been
72    detected in blood[16-19]. Driver mutations have similarly been detected in a high proportion of
73    endometrial crypts[11]. The extent of this phenomenon in the colon, an organ with a high
74    cancer incidence, has not been investigated.
75
76    The colonic epithelial lining is a contiguous cell sheet organised into ~15,000,000 glandular
77    units, known as crypts, oriented perpendicular to the luminal surface and composed of ~2,000
78    cells[20]. Towards the base of each crypt resides a small number of stem cells ancestral to the
79    maturing and differentiated cells in the crypt[21]. These stem cells stochastically replace one
80    another through a process of neutral drift[22,23] such that all stem cells, and thus all cells, in a
81    crypt derive from a single ancestor stem cell that existed in recent years[24-27]. The somatic
82    mutations that were present in this ancestor are thus found in all ~2,000 descendant cells and
83    can be revealed by DNA sequencing of an individual crypt. Following acquisition of the
84    requisite numbers and combinations of driver mutations, these stem cells are also thought to
85    be the cells of origin of colorectal cancers[28]. To characterise the earliest stages of colorectal
86    carcinogenesis, somatic mutation burdens, mutational signatures and the frequency of driver
87    mutations in normal colorectal epithelium were explored by sequencing individual colorectal
88    crypts.
89
90    **Results**
91    **Somatic mutations and mutational signatures**
92    2,035 individual colonic crypts from 42 individuals aged 11 to 78 were isolated using laser
93    capture microdissection and sequenced using a modified library-making protocol developed
94    for small amounts of input DNA. The samples were from seven transplant organ donors, 34
95    individuals biopsied to investigate potential colorectal disease and an autopsy of a subject
96    with oesophageal cancer. In total, 15 had colorectal cancer and 27 showed no evidence of
97    colorectal disease (Supplementary Table). Samples from all individuals in this study are
98    referred to as "normal" crypts as when a cancer was present only biopsies distant from the
99    lesion were used. The distribution of the variant allele fractions of mutations from whole
100   genome sequencing of 571 individual crypts indicated that the large majority of crypts were
101   predominantly clonal cell populations derived from a single ancestral stem cell (Extended
102   Data Fig. 1d). There was substantial variation in mutation burdens between individual crypts,

103    ranging from 1,508 to 15,329 for individuals in their sixties, which was not obviously
104    attributable to technical factors. To explore the biological basis of this variation we extracted
105    mutational signatures and estimated the contribution of each to the mutation burden of each
106    crypt (Methods, Supplementary Results).
107
108    Nine single base substitution (SBS), six doublet base substitution (DBS), and five small indel
109    (ID) mutational signatures were found. Of these, 14 closely matched (Methods) a known
110    reference signature (SBS1, SBS2, SBS5, SBS13, SBS18, DBS2, DBS4, DBS6, DBS8,
111    DBS9, DBS11, ID1, ID2, and ID5, nomenclature as in Alexandrov et al[1]) and six did not
112    (SBSA, SBSB, SBSC, SBSD, IDA, and IDB) (Fig. 1, Extended Data Fig. 2-4). Thus, new
113    mutational signatures were extracted despite extensive prior analysis of cancers, perhaps due
114    to masking by the comparative complexity of signature mixtures present in cancer genomes.
115
116    **Ubiquitous mutational signatures**
117    11 signatures (three SBS, five DBS and three ID) were found in >85% of crypts and are here
118    termed "ubiquitous". All have been previously described[1]. SBS1 is characterised by C>T
119    substitutions at NCG trinucleotides (the mutated base is underlined) and is likely due to
120    deamination of 5-methylcytosine. Its mutation load correlated linearly with age (Fig. 2).
121    There was, however, variation in SBS1 mutation burdens between crypts from the same
122    individual. This was due, in part, to different SBS1 mutation rates in different colonic sectors,
123    with 16.8 mutations per year (95% CI 15.2-18.3) in the right (ascending and caecum), 16.1
124    (95% CI 14.4-17.5) in the transverse, and 12.8 (95% CI 11.1-14.4) in the left colon
125    (descending and sigmoid). The SBS1 mutation rate in the terminal part of the small bowel,
126    the ileum, was 12.7 (95% CI 10.6-14.9) (Supplementary Results).
127
128    SBS5 is a relatively flat, featureless signature of unknown cause and SBS18 is predominantly
129    characterised by C>A mutations, which may be due to DNA damage by reactive oxygen
130    species[29,30]. The mutation burdens of these signatures also showed positive correlations with
131    age, with the same ordering of sector differences as SBS1. Even after taking anatomical
132    location and age into account, differences in mutation burden remained between different
133    crypts, notably for SBS18, indicating that additional factors influence mutation rates in
134    normal cells (Fig. 2, Extended Data Fig. 9, Extended Data Fig. 6al).
135
136    DBS2, DBS4, DBS6, DBS9, and DBS11 were tightly correlated in all colonic crypts. A
137    composite spectrum of DBS2 and DBS4 is also present in normal mouse cells and, in human
138    cancers, both correlate with age of diagnosis confirming that a substantial proportion of their
139    mutations are generated in normal cells[1]. ID1, ID2, and ID5, which are predominantly
140    characterised by insertions and deletions of a single T and may be the consequence of
141    slippage during DNA replication, all accumulated linearly with age with the same order of
142    sector differences as SBS1, SBS5, and SBS18 (Supplementary Results, Extended Data Fig.
143    5).
144
145    The correlations of mutation burden with age indicate that the mutational processes
146    underlying these ubiquitous mutational signatures operate continuously throughout life, in all
147    individuals and in all colorectal stem cells at similar rates. However, the results also suggest
148    that differences in physiology and/or microenvironment (and potentially age of the most
149    recent common ancestor of crypts[27]) between cells in different sectors of the colon cause
150    measurable differences in somatic mutation rates.
151
152    **Sporadic mutational signatures**

153   Nine signatures (six SBS, one DBS and two ID) were present only in a subset of individuals
154   and/or a subset of crypts and are termed "sporadic". All were novel, except for SBS2, SBS13
155   and DBS8. SBS2 and SBS13 are predominantly characterised by C>T and C>G mutations at
156   TCN, are likely due to activity of APOBEC cytidine deaminases and usually occur
157   together[31,32]. SBS2 and SBS13 were clearly observed in one colonic crypt from one
158   individual and one ileal crypt from another, but smaller contributions may be present in
159   additional crypts (Extended Data Fig. 9). To our knowledge, this is the first reported evidence
160   that APOBEC DNA-editing of the human genome occurs in normal cells *in vivo*. However,
161   in the colon at least, it is restricted to a small subset of cells. The factors that initiate it are
162   unknown, although viral entry, retrotransposon transposition and local inflammation have
163   been proposed in other contexts[33]. The wider sequence context of these mutations in normal
164   colon suggests that APOBEC3A is the major contributing enzyme[34].
165
166   Four SBS signatures that do not match the reference set, SBSA-D, were found in normal
167   colorectal cells (SBSA has recently been reported in an oral squamous carcinoma[35]). SBSA is
168   characterised predominantly by T>C mutations at ATA, ATT, and TTT, and T>G mutations
169   at TTT. Its mutation burden correlated closely with that of IDA, which is characterised by
170   single T deletions in short runs of Ts (with a mode of four), suggesting that they are due to
171   the same underlying mutational process. SBSA exhibited a highly variable mutational
172   burden, being present in 29/42 individuals studied, often in just a subset of crypts, and
173   showed evidence of spatial clustering in the colon, with crypts from the same biopsy carrying
174   the signature even though the mutations themselves were not shared (Supplementary Results,
175   Extended Data Fig. 9). 2.5-fold more T>C mutations occurred when the T was on the
176   transcribed than on the untranscribed strand. Transcriptional strand bias is often due to
177   transcription coupled nucleotide excision repair (TC-NER) acting on DNA damaged by
178   exogenous exposures causing covalently bound bulky adducts, but can also be caused by
179   transcription coupled DNA damage[36]. Assuming either is the case, damage to adenine
180   underlies SBSA. To investigate the timing of SBSA mutation generation, phylogenetic trees
181   of mutations were constructed and the mutational signatures in each branch established (Fig.
182   3, Extended Data Fig. 6). SBSA was confined to early branches of these phylogenies (when
183   these were available for analysis). (Fig. 3b, Extended Data Fig. 6 f, h, z, aa, am, ao). Using
184   the number of SBS1 mutations as indicators of real time, the mutational process underlying
185   SBSA appears to be active usually before 10 years of age. The initiating event for this
186   relatively frequent mutational process is unknown, but the results suggest an extrinsic, locally
187   acting and patchily distributed mutagenic insult occurring during childhood.
188
189   SBSB was predominantly characterised by C>T substitutions at ACA, T>A at CTN, and T>G
190   at GTG and was present in subsets of crypts from four individuals (Fig. 3c, Extended Data
191   Fig. 6). In the two individuals in whom it could be timed (Extended Data Fig. 6 aa, ai), it
192   appeared – as with SBSA – to be active in the first decade of life. SBSB correlated with
193   DBS8 and IDB (Fig. 3c, Extended Data Fig. 9), suggesting that they are caused by the same
194   underlying mutational process. DBS8 is composed of AC>CA and AC>CT mutations and
195   has previously been reported in rare hypermutated cancers with no obvious cause[1]. IDB is
196   dominated by deletion of a single T with no other Ts surrounding it. The mutational process
197   underlying this signature is unknown.
198
199   SBSC is characterised predominantly by one C>T mutation in CC dinucleotides. It is of
200   unknown aetiology and primarily affects three crypts from the left colon of one individual
201   with an unremarkable history (Extended Data Fig. 9).
202

203    All crypts from a 66 year-old man carried many thousands of mutations of SBSD,
204    characterised predominantly by T>A substitutions with a transcriptional strand bias
205    compatible with damage to adenine. This individual had been treated with multiple
206    chemotherapeutic agents (including cyclophosphamide, doxorubicin, vincristine,
207    prednisolone, chlorambucil, bleomycin and etoposide) for lymphoma and subsequently
208    developed caecal adenocarcinoma. SBSD resembles SBS25, (cosine similarity 0.9),
209    previously found in Hodgkin lymphoma cell lines from two chemotherapy-treated
210    patients[31,37]. To our knowledge this is the first time that the mutational consequences of
211    chemotherapy have been demonstrated in normal human cells *in vivo*. The mutation burden
212    in this individual's colorectal epithelium was 3-5 fold higher than expected for his age, thus
213    by extrapolation equivalent to that of a 200-300 year-old, and it is plausible that other tissues
214    have been similarly affected.
215
216    **Copy number changes and structural variants**
217    Copy number changes and/or structural variants were found in 80 out of 449 (18%) evaluable
218    normal crypts. Five crypts exhibited eight whole chromosome copy number increases which,
219    notably, affected the same three chromosomes: 3, 7 and 9, as well as the X chromosome
220    (Extended Data Fig. 7a). Thus, copy number increases clustered in certain crypts and tended
221    to affect certain chromosomes. No whole chromosome losses were observed. Regions of
222    copy number neutral loss of heterozygosity were observed in 12 crypts, affecting
223    chromosomes 1p, 6p, 7p, 8q, 9q, 10q (twice), 17p, 17q, 18q, 21q and 22q. Five of these copy
224    number changes could be timed and all were estimated to have occurred in adulthood. Two
225    changes that affected the same crypt appeared to be synchronous (Supplementary Results,
226    Extended Data Fig. 7b). Forty-eight large deletions, 18 tandem duplications, four
227    translocations, and two inversions were detected. All were private to a single crypt, except for
228    one deletion which was present in two adjacent crypts sharing few mutations, indicating that
229    the deletion occurred during gestation or early childhood.
230
231    **Driver mutations**
232    Driver mutations are those that confer a selective advantage during cancer evolution and
233    may, but need not, promote neoplasia[38]. To search for driver mutations in normal colon, the
234    whole genome sequences of 571 crypts were supplemented with targeted sequencing of 90
235    known colorectal cancer genes in an additional series of crypts. In total, substitutions in these
236    genes were evaluable in 1,403 crypts and indels in 1,046. Statistical analysis revealed
237    evidence of positive selection on the recessive cancer genes *AXIN2* (three truncating
238    mutations, adjusted q value 0.004) and *STAG2* (two truncating mutations, adjusted q value
239    0.038) indicating that these mutations are likely drivers. Additional likely driver mutations
240    were identified in cancer genes characterised by canonical missense hotspot mutations. Nine
241    hotspot mutations in *PIK3CA* (E542K, R38H), *ERBB2* (R678Q, V842I, T862A), *ERBB3*
242    (R475W, R667L), and *FBXW7* (R505C, R658Q) were observed (Extended Data Fig. 8).
243    Given the specificity of these hotspot mutations, most are likely to be drivers. In addition,
244    heterozygous truncating mutations were found in the recessive cancer genes *ARID2*, *ATM*
245    (two), *ATR*, *BRCA2*, *CDK12* (two), *CDKN1B*, *RNF43* (two), *TBL1XR1*, and *TP53*
246    (Supplementary Table). There was no statistical evidence for selection of truncating
247    mutations in the set of 90 colorectal cancer genes overall. The possibility that some have
248    conferred clonal growth advantage, however, is not excluded. No crypt carried more than one
249    putative driver mutation.
250
251    23 pairs of adjacent crypts shared over 100 SBS1 mutations and thus were likely to have been
252    generated by postnatal crypt fission. Two pairs carried driver mutations (one with an *AXIN2*

253    nonsense mutation and one with *PIK3CA* E542K), although the association of driver
254    mutations with crypt fission is not significant (p=0.17). In one sister crypt the *AXIN2*
255    mutation was rendered homozygous by copy number neutral chromosome 17q LOH,
256    revealing ongoing clonal evolution in normal colon (Fig. 4, Fig. 3b).
257
258    On the conservative assumption that just the *AXIN2* and *STAG2* truncating mutations and the
259    missense hotspot mutations in *PIK3CA*, *ERBB2*, *ERBB3* and *FBXW7* are drivers, ~1% of
260    normal colorectal crypts (~150,000 crypts) in a 50-60 year old carries a driver mutation.
261    Since in the over 70s ~40% of people have an adenoma on colonoscopy[39] and ~5% of people
262    develop colorectal cancer over their lifetime[40] (and some of these may arise from more
263    recently-acquired driver mutations) only an extremely small proportion of these crypt
264    microneoplasms becomes a macroscopically detectable adenoma (< 1/375,000) or carcinoma
265    (< 1/3,000,000) within the following few decades.
266
267    The proportion of normal colorectal cells with a driver mutation (1%) is considerably lower
268    than that observed in normal skin (30%). The lower frequency of drivers in colon may be
269    due, at least in part, to the modular structure of glandular epithelia. The small number of stem
270    cells within a crypt diminishes the probability that a cell with a driver mutation will
271    outcompete its wild-type neighbours. Moreover, even if it does colonise the crypt, a mutant
272    stem cell is entombed in it unless it can overcome the largely unknown forces that govern
273    clonal expansion through crypt fission.
274
275    **Comparisons with colorectal cancer**
276    There are marked differences between the genomes of normal colorectal stem cells and those
277    of colorectal cancers. The total mutation burdens of base substitutions (10,000-20,000) and
278    indels (1,000-2,000) found in most colorectal carcinomas[1] (excluding those with
279    hypermutator phenotypes in which it is usually >10-fold more) is higher than the ~3,000
280    substitutions and 300 indels found in most normal crypts from 50-60 year old individuals.
281    The particularly high base substitution and indel mutation burdens and associated mutational
282    signatures of DNA mismatch repair deficiency and/or polymerase epsilon/delta mutations
283    were not found in any normal colorectal crypts but are present in ~20% colorectal cancers.
284    Equally striking is the difference between the 0-4 structural changes per normal crypt (with
285    the majority having none) and the 10s to 100s per colorectal cancer[41]. In all these respects,
286    the genomes of normal crypts with driver mutations were similar to those of normal crypts
287    without drivers (Extended Data Fig. 9).
288
289    Elevated mutation burdens are, therefore, characteristic of the evolutionary trajectory from
290    normal colorectal cell to cancer cell. The increased base substitution and indel mutation loads
291    in cancers are due to a combination of higher burdens of the ubiquitous mutational signatures
292    found in normal crypts, additional base substitution and indel signatures thus far found
293    exclusively in cancers (confirming previous reports[5,42]) and larger numbers of copy number
294    changes and structural variation. The causes of some of these additional mutational loads are
295    known (for example, defective DNA mismatch repair and polymerase epsilon/delta
296    mutations) but the remainder are uncertain.
297
298    The relative frequencies of mutated cancer genes differ between colorectal
299    adenomas/carcinomas and normal colorectal cells (p=0.003, Supplementary Results). In
300    colorectal cancer, mutations in *APC*, *KRAS* and *TP53* are common[43], accounting for 56% of
301    base substitution and indel drivers (Supplementary Methods) but are comparatively rare
302    among normal crypts with driver mutations (1/14). By contrast, mutations in, for example,

303    *ERBB2* and *ERBB3* are relatively common in normal crypts with drivers (5/14) but rare in
304    colorectal cancer (7/631). The results suggest that mutations in *APC*, *KRAS* and *TP53* confer
305    higher likelihoods of conversion to adenoma and carcinoma than mutations in *ERBB2* and
306    *ERBB3* whereas mutations in *ERBB2* and *ERBB3* confer higher likelihoods of stem cells
307    colonising crypts than *APC*, *KRAS* and *TP53*. Nevertheless, previous reports suggest that
308    1:3,500 epithelial cells, and therefore >4,000 crypts per colon, bear *KRAS* G12D[44], and so
309    even these have a low probability of progression.
310
311

312    **Discussion**
313    This study has characterised all classes of somatic mutation in normal colorectal epithelial
314    stem cells. A substantial repertoire of base substitution and indel mutational processes is
315    operative, some ubiquitous and some sporadic, together with relatively infrequent copy
316    number changes and genome rearrangements. APOBEC DNA-editing occurs in normal
317    colon, albeit only in rare cells. Many signatures, however, are of unknown aetiology and
318    some appear to be acquired early in life. The presence of five times the age-standard mutation
319    load in all colorectal cells, and potentially many other tissues, in an individual who had
320    undergone chemotherapy provides new insight into the impact of such exposures and raises
321    questions pertaining to its relationship with chemotherapy's relatively modest impact on
322    cancer risk[45].
323

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

336    The proportion of normal colorectal epithelial cells with driver mutations is, however,
337    substantially lower than that of other normal tissues so far studied, notably skin[10] and
338    endometrium[11]. Colorectal epithelium is constituted of crypts, modular units which may
339    themselves constrain clonal expansion, and this architecture may contribute to such
340    differences with skin. The reason for the difference with endometrium, which is also
341    glandular, remains to be explored.
342

343    Fundamental questions are being addressed with respect to differences in cancer incidence
344    rates between tissues. The somatic mutation burden in colon and ileum is similar despite the
345    substantially higher cancer incidence rate in colon (as previously noted[4]) and therefore does
346    not appear to account for this difference. Whether the total burden of microneoplastic change
347    across the colon and in other tissues more closely correlates with these differences is yet to be
348    determined.
349

350    Finally, this study provides a reference perspective on the mutational signatures and driver
351    mutations in normal colon against which disease states of inflammatory, genetic, neoplastic,
352    degenerative and other aetiologies can be compared. Similar surveys conducted across the

353 range of normal cell types will inform on the universal process of somatic evolution in the
354 human body in health and disease.
355
356
357
## ACKNOWLEDGEMENTS

367
368
369
## AUTHOR CONTRIBUTIONS

371 MRS and HLS designed the study and wrote the manuscript with contributions from all the
372 authors. KSP, NC, MZ, RCF, NG, FT, AN, MG, and LM recruited patients and obtained
373 samples. PE, RO, HLS, and LM devised the protocol to laser capture microdissect and
374 sequence colonic crypts. HLS prepared sections, microdissected, and lysed colonic crypts. PR
375 contributed to laser capture microdissection. PE and CA made libraries. HLS performed most
376 of the data curation and statistical analysis. MAS devised filters for substitution calling. JW
377 performed in-house NMF signature extraction. TC contributed to statistical analyses. LON
378 provided technical assistance. PJC and IM oversaw statistical analyses. MRS supervised the
379 study.
380
381
## REFERENCES

383
384 1. Alexandrov, L.B. et al. The repertoire of mutational signatures in human cancer.
385    Preprint at: https://www.biorxiv.org/content/early/2018/05/15/322859 (2018).
386 2. Sabarinathan, R. et al. The whole genome panorama of cancer drivers. Preprint at:
387    https://www.biorxiv.org/content/early/2017/12/23/190330 (2017).
388 3. Fearon E.R. & Vogelstein B. A genetic model for colorectal tumorigenesis. *Cell* **61**,
389    759-767 (1990).
390 4. Blokzijl, F. et al. Tissue-specific mutation accumulation in human adult stem cells
391    during life. *Nature* **538**, 260–264 (2016).
392 5. Roerink S.F. et al., Intra-tumour diversification in colorectal cancer at the single cell
393    level. *Nature* **556,** 457-462 (2018).
394 6. Welch, J. S. et al. The origin and evolution of mutations in acute myeloid leukemia.
395    *Cell* **150**, 264–278 (2012).
396 7. Bae, T. et al. Different mutational rates and mechanisms in human cells at
397    pregastrulation and neurogenesis. *Science* **359**, 550–555 (2018).
398 8. Behjati, S. et al. Genome sequencing of normal cells reveals developmental lineages
399    and mutational processes. *Nature* **513**, 422–425 (2014).
400 9. Lee-Six H. et al., Population dynamics of normal human blood inferred from somatic
401    mutations. *Nature* https://doi.org/10.1038/s41586-018-0497-0 (2018).

402  10. Martincorena, I. et al. Tumor evolution. High burden and pervasive positive
403      selection of somatic mutations in normal human skin. *Science* **348**, 880–886
404      (2015).
405  11. Suda, K. et al. Clonal Expansion and Diversification of Cancer-Associated
406      Mutations in Endometriosis and Normal Endometrium. *Cell Rep.* **24**, 1777-1789
407      (2018).
408  12. Xu, X. et al. Single-cell exome sequencing reveals single-nucleotide mutation
409      characteristics of a kidney tumor. *Cell* **148**, 886-895 (2012).
410  13. Lodato M.A. et al. Somatic mutation in single human neurons tracks developmental
411      and transcriptional history. *Science* **350**, 94-98 (2015).
412  14. Lodato, M. A. et al. Aging and neurodegeneration are associated with increased
413      mutations in single human neurons. *Science* **359**, 555–559 (2018).
414  15. Hoang, M.L. et al. Genome-wide quantification of rare somatic mutations in normal
415      human tissues using massively parallel sequencing. *PNAS* **113**, 9846-9851 (2016).
416  16. Jaiswal, S. et al. Age-related clonal hematopoiesis associated with adverse outcomes.
417      *N. Engl. J. Med.* **371**, 2488–2498 (2014).
418  17. Xie, M. et al. Age-related mutations associated with clonal hematopoietic expansion
419      and malignancies. *Nat. Med.* **20**, 1472–1478 (2014).
420  18. McKerrell, T. et al. Leukemia-associated somatic mutations drive distinct patterns of
421      age-related clonal hemopoiesis. *Cell Rep.* **10**, 1239–1245 (2015).
422  19. Genovese, G. et al. Clonal hematopoiesis and blood-cancer risk inferred from blood
423      DNA sequence. *N. Engl. J. Med.* **371**, 2477-2487 (2015).
424  20. Potten, C.S. et al. Measurement of in vivo proliferation in human colorectal mucosa
425      using bromodeoxyuridine. *Gut.* **33**, 71-78 (1992).
426  21. Cheng, H. & Leblond, C.P. Origin, differentiation and renewal of the four main
427      epithelial cell types in the mouse small intestine. V. Unitarian Theory of the origin of
428      the four epithelial cell types. *Am J Anat.* **141**, 537-561 (1974).
429  22. Lopez-Garcia, C. et al. Intestinal stem cell replacement follows a pattern of neutral
430      drift. *Science* **330**, 822-825 (2010).
431  23. Snippert, H.J. et al. Intestinal crypt homeostasis results from neutral competition
432      between symmetrically dividing Lgr5 stem cells. *Cell* **143**, 134-144 (2010).
433  24. Griffiths D.F. et al. Demonstration of somatic mutation and colonic crypt clonality by
434      X-linked enzyme histochemistry. *Nature* **333**, 461-463 (1988).
435  25. Winton, D.J., and Ponder, B.A. Stem-cell organization in mouse small intestine. *Proc.*
436      *Biol. Sci.* **241**, 13-18 (1990).
437  26. Kozar, S. et al. Continuous clonal labeling reveals small numbers of functional stem
438      cells in intestinal crypts and adenomas. *Cell Stem Cell* **13**, 626-633 (2013).
439  27. Nicholson, A. et al. Fixation and spread of somatic mutations in adult human colonic
440      epithelium. *Cell Stem Cell* **22**, 909-918 (2018).
441  28. Barker, N. et al. Crypt stem cells as the cells-of-origin of intestinal cancer. *Nature*
442      **457**, 608-611 (2009).
443  29. Rouhani, F.J. et al. Mutational history of a human cell lineage form somatic to
444      induced pluripotent stem cells. *PLoS Genet.* **12**, e1005932 (2016).
445  30. Viel, A. et al. A specific mutational signature associated with DNA 8-Oxoguanine
446      persistence in MUTYH-defective colorectal cancer. *EBioMedicine* **20**, 39-49 (2017).
447  31. Alexandrov, L.B. et al. Signatures of mutational processes in human cancer. *Nature*
448      **500**, 415-421 (2013).
449  32. Nik-Zainal, S. et al. Mutational processes molding the genomes of 21 breast cancers.
450      *Cell* **149**, 979-993 (2012).

451     33. Vieira, V.C. & Soares, M.A. The role of cytidine deaminases on innate immune
452         responses against human viral infections. *Biomed Res Int*, **683095** (2013).
453     34. Chan, K. et al. An APOBEC3A hypermutation signature is distinguishable from the
454         signature of background mutagenesis by APOBEC3B in human cancers. *Nat Genet.*
455         **47**, 1067-1072 (2015).
456     35. Boot, A. et al. Mutational signature analysis of Asian OSCCs reveals novel
457         mutational signature with exceptional sequence context specificity. Preprint at:
458         https://www.biorxiv.org/content/early/2018/07/19/368753.1 (2018).
459     36. Haradhvala, N. J. et al. Mutational Strand Asymmetries in Cancer Genomes Reveal
460         Mechanisms of DNA Damage and Repair. *Cell* **164**, 538-549 (2016).
461     37. Wolf, J. et al. Peripheral blood mononuclear cells of a patient with advanced
462         Hodgkin's lymphoma give rise to permanently growing Hodgkin-Reed Sternberg
463         cells. *Blood* **87**, 3418-3428 (1996).
464     38. Stratton, M.R., Campbell, P.J., & Futreal, P.A. The cancer genome. *Nature* **458**, 719-
465         724 (2009).
466     39. Corley, D.A., et al. Variation of adenoma prevalence by age, sex, race, and colon
467         location in a large population: implications for screening and quality programs. *Clin.*
468         *Grastroenterol. Hepatol.* **11**, 172-180 (2013).
469     40. Cancer Research UK, Bowel Cancer Incidence Statistics,
470         https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-
471         cancer-type/bowel-cancer/incidence#heading-Seven (Accessed August 2018).
472     41. Li, Y. et al. Patterns of structural variation in human cancer. Preprint at:
473         https://www.biorxiv.org/content/early/2017/08/27/181339 (2017).
474     42. Lugli, N. et al. Enhanced Rate of Acquisition of Point Mutations in Mouse Intestinal
475         Adenomas Compared to Normal Tissue. *Cell Reports* **19**, 2185-2192 (2017).
476     43. The Cancer Genome Atlas Network, Comprehensive molecular characterization of
477         human colon and rectal cancer. *Nature* **487**, 330-337 (2012)
478     44. Parsons, B.L. et al. ACB-PCR quantification of K-RAS codon 12 GAT and GTT
479         mutant fraction in colon tumor and non-tumor tissue. *Cancer Invest.* **28**, 364-375
480         (2010).
481     45. Travis, L.B. Therapy-associated solid tumors, *Acta Oncologica*, **41**, 323-333 (2002).

482
483
484
485     **FIGURE LEGENDS**
486
487     **Figure 1. Mutational signatures present in normal colon. a**, an example SBS, DBS, and
488     ID signature showing the categories into which mutations are divided. Later figures are
489     shown with the same categories, ordering, and colour scheme. **b**, the complement of
490     signatures discovered in normal colonic epithelium. Known signatures are labelled according
491     to their nomenclature in PCAWG, while novel signatures are labelled with letters. SBS,
492     single base substitution; DBS, doublet base substitution; ID, small insertion or deletion.
493
494     **Figure 2. Mutation burden *versus* age for every signature.** For every signature, the median
495     (horizontal bar) and range (vertical bar) in mutation burden for all the crypts from each
496     individual are shown. Each individual is coloured differently. See Supplementary Results for
497     plots showing every crypt.
498
499     **Figure 3. Crypt phylogenies.** For four selected individuals (**a-d**), the phylogeny is shown
500     three times: on top, with branch lengths proportional to the number of single base

501 substitutions; in the middle, with branch lengths proportional to the number of doublet base
502 substitutions; on the bottom, with branch lengths proportional to the number of small
503 insertions and deletions. Scale bars are shown on the right-hand side. A stacked barplot of the
504 mutational signatures that contribute to each branch is superimposed onto every branch.
505 Please note that the ordering of signatures along a given branch is just for visualisation
506 purposes: we cannot distinguish the timing of different signatures along a branch. "X0"
507 indicates mutations that could not confidently be assigned to any signature. The phylogenies
508 for all individuals are shown in Extended Data Fig. 6. (**a**) a phylogeny dominated by
509 ubiquitous and known signatures. A *PIK3CA* mutation is shared by two crypts. (**b**) a
510 phylogeny with a strong contribution of SBSA and IDA, as well as an *AXIN2* mutation (the
511 same as in Fig. 4). (**c**) a phylogeny with SBSB, DBS8, and IDB. (**d**) the phylogeny of the
512 individual exposed to chemotherapy, showing a strong contribution of SBSD.
513
514 **Figure 4. An *AXIN2* driver mutation in normal colon.** (**a**) a section (after dissection) in
515 which an inactivating *AXIN2* mutation was found. Red dots represent crypts with the *AXIN2*
516 mutation. Blue dots represent crypts that could be assessed and were found not to have the
517 mutation. Crypts without dots failed sequencing and could not be assessed. (**b**) the two crypts
518 with the *AXIN2* mutations prior to dissection did not appear different to other crypts. (**c**) copy
519 neutral loss of heterozygosity (CNN-LOH) of one of the crypts over the *AXIN2* locus. The
520 copy number state (y axis) for every chromosome is shown, with one allele coloured red and
521 the other green. (**d**) Jbrowse image of reads supporting the *AXIN2* mutations in each of the
522 crypts. The mutation is coloured red. 25 out of 29 reads support the mutation in the crypt that
523 has CNN-LOH; the four reads that do not are presumably the result of stromal contamination.
524
525 **Figure 5. Comparison of the mutational signatures and driver landscape of normal
526 crypts and colorectal adenocarcinomas.** (**a**) a comparison of the burden of mutations due to
527 every mutational signature found in either group. For each signature, the (mutation burden+1)
528 of every sample is shown on the y axis on a log scale. Normal colon and cancer samples are
529 ordered within their groups. Colorectal adenocarcinoma signature attributions and burden are
530 from Alexandrov et al.[1]. (**b-c**) the frequency of driver mutations in normal colon and
531 colorectal cancer. The frequency of driver mutations is derived using data from The Cancer
532 Genome Atlas Network[43] (Supplementary Methods). (**b**) the proportion of crypts or cancers
533 with driver mutations in each gene found in either of the two groups. (**c**) the proportion of
534 driver mutations in each gene in normal and cancer.
535
536
537
538 **EXTENDED FIGURE LEGENDS**
539
540 **Extended Data Figure 1. Laser capture microdissection of crypts.** (**a**) a representative
541 image of a section of colonic tissue, with a magnified inset showing the section before and
542 after dissection of a crypt. (**b-c**), the coverage of crypts that underwent whole genome (**b**) and
543 targeted (**c**) sequencing. (**d-e**), their respective VAF (which is half of the clonal fraction).
544
545 **Extended Data Figure 2. HDP signature extraction results.** Results of signature
546 extraction using an HDP with conditioning on signatures known to be active in colorectal
547 cancer. For each signature, the extracted signature and the profile of a sample that has a
548 strong contribution of that signature are shown. Signatures are presented as in Fig. 2. The
549 HDP extraction was followed by deconvolution by Expectation Maximisation (Methods,

550  Extended Data Fig. 3) to produce the version of signatures presented in the main text. HDP,
551  Hierarchical Dirichlet Process.
552
553  **Extended Data Figure 3. Expectation maximisation decomposition of HDP signatures.**
554  Three signatures were decomposed. For each panel, the original HDP version in shown on
555  the top left, the PCAWG signatures that are deemed to contribute at least 10% of mutations to
556  it on the right, and the reconstituted signature built by combining the PCAWG signatures on
557  the bottom left. The cosine similarity of the reconstituted signature to the original is shown in
558  the title to the reconstituted signature plot. HDP, Hierarchical Dirichlet Process; PCAWG,
559  Pan Cancer Analysis of Whole Genomes.
560
561  **Extended Data Figure 4. Validation of single base substitution signatures.** Other methods
562  of signature extraction were run to test the robustness of signature decomposition. **a**, HDP
563  without pre-conditioning on PCAWG. **b**, In-house NNMF without pre-conditioning on
564  PCAWG. **c**, NNMF implemented by the MutationalPatterns R package (Methods). HDP,
565  Hierarchical Dirichlet Process; PCAWG, Pan Cancer Analysis of Whole Genomes; NNMF,
566  Non-Negative Matrix Factorisation.
567
568  **Extended Data Figure 5. Linear modelling of signature accumulation.** For signatures that
569  appeared to show a linear accumulation with age, the mutation rate per site was determined
570  using mixed models, with age and site as fixed effects, and individual as a random effect.
571  Confidence intervals were determined by bootstrapping.
572
573  **Extended Data Figure 6. Crypt phylogenies.** For every individual, the phylogeny of crypts
574  is shown three times: on top, with branch lengths proportional to the number of single base
575  substitutions; in the middle, with branch lengths proportional to the number of doublet base
576  substitutions; on the bottom, with branch lengths proportional to the number of small
577  insertions and deletions. Scale bars are shown on the right-hand side. A stacked barplot of the
578  mutational signatures that contribute to each branch is overlaid over every branch. "X0"
579  indicates mutations that could not confidently be assigned to any signature. Please note that
580  the ordering of signatures along a given branch is just for visualisation purposes: we cannot
581  distinguish the timing of different signatures along a branch.
582
583  **Extended Data Figure 7. Copy number changes in normal colon.** (**a**) whole chromosome
584  amplifications in five crypts. The copy number state (y axis) for each allele, one coloured red,
585  and one coloured green, is shown. Chromosomes are labelled along the top of the graph. (**b**)
586  timing of copy number changes throughout life. Vertical bars represent 95% confidence
587  intervals determined by bootstrapping.
588
589  **Extended Data Figure 8. Gain of function driver mutations in normal colon.** Putative
590  driver missense mutations in oncogene hotspots. The number of substitutions catalogued in
591  COSMIC are shown on the y axis at each position along the gene, with the mutations
592  observed in our cohort highlighted.
593
594  **Extended Data Figure 9. Occurrence matrix of signatures and driver mutations in**
595  **crypts.** For all crypts that were whole genome sequenced to sufficient depth and for crypts
596  that underwent targeted sequencing and in which driver mutations were found, the signatures
597  and driver mutations are shown. Each vertical column represents a crypt. The individual to
598  which each crypt belongs is indicated by alternating colours in the top bar. The site to which
599  each crypt belongs is shown underneath. The contribution of each signature to each crypt;

600  thus the crypt with the largest contribution of a given signature is coloured purple, and the
601  crypt with the smallest contribution is coloured white. Crypts in which the signatures could
602  not be assessed, either because they underwent targeted sequencing or the coverage was poor,
603  are coloured grey. Driver mutations, including heterozygous mutations in tumour suppressor
604  genes, are indicated by a black bar.
605
606
607
608  **SUPPLEMENTARY METHODS**
609
610  **Human tissues**
611  We obtained healthy colonic biopsies from four cohorts. The first represents seven deceased
612  organ donors ranging in age from 36 to 67, from whom colonic and small intestinal biopsies
613  were taken at the time of organ donation (REC 15/EE/0152). The second represents
614  individuals aged 60 to 72 who were having a colonoscopy following a positive faecal occult
615  blood test as part of the Bowel Cancer Screening Programme (Ethical approval 08-H0308-
616  13); we selected 16 who were not found to have either an adenoma or a carcinoma on
617  colonoscopy, and 15 who were found to have a colorectal carcinoma (the normal biopsies
618  that we use were distant from these lesions). The third cohort represents three paediatric
619  patients who underwent routine colonoscopy to exclude inflammatory bowel disease and who
620  were found to have a completely normal intestinal mucosa macroscopically and histologically
621  (REC 12/EE/0482). The final cohort included one 78 year-old gentleman with oesophageal
622  cancer who underwent a warm autopsy (REC 13/EE/0043). All samples were obtained with
623  informed consent and studies approved by East of England Research Ethics Committees.
624
625  **Laser capture microdissection of colonic crypts**
626  Fresh frozen biopsies were embedded in optimal cutting temperature (OCT) compound. 30
627  micrometre sections were fixed in methanol for five minutes, washed three times with
628  phosphate-buffered saline, and stained with Gill's haematoxylin for 20 seconds. Crypts were
629  isolated by laser capture microdissection, and collected in separate wells of a 96-well plate.
630  They were lysed using the Arcturus PicoPure Kit (Applied Biosystems) according to the
631  manufacturer's instructions. DNA library prep then proceeded without clean-up or
632  quantification.
633
634  **Library preparation**
635  Two library preparation methods were used for laser capture microdissected (LCM) material:
636  in initial experiments sonication was used to fragment DNA, and later, an enzymatic
637  fragmentation method was implemented as it could make libraries from even lower input.
638  Comparison of the two methods showed no difference in mutation calls once post-processing
639  filters (described below) had been implemented. All samples in this study were processed
640  using an Agilent Bravo Workstation (Option B; Agilent Technologies).
641
642  For sonication libraries, LCM lysate (20 µl) was mixed with 100 µl TE buffer (Ambion; 10
643  mM Tris-HCl, 1 mM EDTA) and DNA was fragmented using focused acoustics (Covaris
644  LE220; Covaris, Inc.). Fragmented DNA was mixed with 80 µl Ampure XP beads (Beckman
645  Coulter). Following a 5 min binding reaction and magnetic bead separation, genomic DNA
646  was washed twice with 75% ethanol. Beads were resuspended in 20 µl nuclease-free water
647  (Ambion) and processed immediately for DNA library construction. Each sample (20 µl) was
648  mixed with 2.8 µl of NEBNext Ultra II End Prep Reaction Buffer, 1.25 µl of NEBNext Ultra
649  II End Prep Enzyme Mix (New England BioLabs) and incubated on a thermal cycler for 30

650 min at 20°C then 30 min at 65°C.  Following DNA fragmentation and A-tailing, each sample
651 was incubated for 20 min at 20°C with a mixture of 30 μl ligation mix and 1 μl ligation
652 enhancer (New England BioLabs), 0.9 μl nuclease-free water (Ambion) and 0.1 μl duplexed
653 adapters (100 uM; 5'-ACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3', 5'-phos-
654 GATCGGAAGAGCGGTTCAGCAGGAATGCCGAG-3').  Adapter-ligated libraries were
655 purified using Ampure XP beads by addition of 65 μl Ampure XP solution (Beckman
656 Coulter) and 65 μl TE buffer (Ambion).  Following elution and bead separation, DNA
657 libraries (21.5 μl) were amplified by PCR by addition of 25 μl KAPA HiFi HotStart
658 ReadyMix (KAPA Biosystems), 1 μl PE1.0 primer (100 μM; 5'-
659 AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGA
660 TC*T-3')         and         2.5         μl         iPCR-Tag         (40         μM;         5'-
661 CAAGCAGAAGACGGCATACGAGATXGAGATCGGTCTCGGCATTCCTGCTGAACC
662 GCTCTTCCGATC-3') where 'X' represents one of 96 unique 8-base indexes The sample
663 was then mixed and thermal cycled as follows: 98 °C for 5 min, then 12 cycles of 98 °C for
664 30 s, 65°C for 30 s, 72 °C for 1 min and finally 72 °C for 5 min. Amplified libraries were
665 purified using a 0.7:1 volumetric ratio of Ampure Beads (Beckman Coulter) to PCR product
666 and eluted into 25 μl of nuclease-free water (Ambion).  DNA libraries were adjusted to 2.4
667 nM and sequenced on the HiSeq X platform (illumina) according to the manufacturer's
668 instructions     with     the     exception     that     we     used     iPCRtagseq     (5'-
669 AAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTC-3') to read the library index.
670
671 For enzymatic fragmentation, LCM lysate (20 ul) was mixed with 50 ul Ampure XP beads
672 (Beckman Coulter) and 50 μl TE buffer (Ambion; 10 mM Tris-HCl, 1 mM EDTA) at room
673 temperature.  Following a 5 min binding reaction and magnetic bead separation, genomic
674 DNA was washed twice with 75% ethanol.  Beads were resuspended in 26 μl TE buffer and
675 the bead/genomic DNA slurry was processed immediately for DNA library construction.
676 Each sample (26 μl) was mixed with 7 μl of 5X Ultra II FS buffer, 2 μl of Ultra II FS enzyme
677 (New England BioLabs) and incubated on a thermal cycler for 12 min at 37°C then 30 min at
678 65°C.  Following DNA fragmentation and A-tailing, each sample was incubated for 20 min at
679 20°C with a mixture of 30 μl ligation mix and 1 μl ligation enhancer (New England
680 BioLabs), 0.9 μl nuclease-free water (Ambion) and 0.1 μl duplexed adapters (100 uM; 5'-
681 ACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3',                               5'-phos-
682 GATCGGAAGAGCGGTTCAGCAGGAATGCCGAG-3').  Adapter-ligated libraries were
683 purified using Ampure XP beads by addition of 65 μl Ampure XP solution (Beckman
684 Coulter) and 65 μl TE buffer (Ambion).  Following elution and bead separation, DNA
685 libraries (21.5 μl) were amplified by PCR by addition of 25 μl KAPA HiFi HotStart
686 ReadyMix (KAPA Biosystems), 1 μl PE1.0 primer (100 μM; 5'-
687 AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGA
688 TC*T-3')         and         2.5         μl         iPCR-Tag         (40         μM;         5'-
689 CAAGCAGAAGACGGCATACGAGATXGAGATCGGTCTCGGCATTCCTGCTGAACC
690 GCTCTTCCGATC-3') where 'X' represents one of 96 unique 8-base indexes The sample
691 was then mixed and thermal cycled as follows: 98 °C for 5 min, then 12 cycles of 98 °C for
692 30 s, 65°C for 30 s, 72 °C for 1 min and finally 72 °C for 5 min. Amplified libraries were
693 purified using a 0.7:1 volumetric ratio of Ampure Beads (Beckman Coulter) to PCR product
694 and eluted into 25 μl of nuclease-free water (Ambion).  DNA libraries were adjusted to 2.4
695 nM and sequenced on the HiSeq X platform (Illumina) according to the manufacturer's
696 instructions     with     the     exception     that     we     used     iPCRtagseq     (5'-
697 AAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTC-3') to read the library index.
698
699 **Whole genome sequencing**

700 We generated paired end sequencing reads (150bp) using Illumina XTEN® machines
701 resulting in ~15x coverage per sample. Sequences were aligned to the human reference
702 genome (NCBI build37) using BWA-MEM.
703
704 **Targeted sequencing**
705 A 2.3 MB capture panel was designed in-house to pull down genes that are known or
706 suspected to play a role in neoplasia. We performed custom RNA bait design following the
707 manufacturer's guidelines (SureSelect, Agilent). Samples were multiplexed on flow cells and
708 subjected to paired end sequencing (75-bp reads) using Illumina HiSeq2000 machines. One
709 96-well plate of samples was sequenced on each lane, but as tissue recovery was variable, a
710 range of coverage was achieved. Sequences were aligned to the human reference genome
711 (NCBI build37) using BWA-align.
712
713 **Data Availability**
714 Whole genome and targeted sequencing data are deposited in the European Genome
715 Phenome Archive (EGA). sequencing data have been deposited with EGA accession
716 EGAD00001004192, EGAD00001004192, and EGAD00001004193.
717
718 **Code Availability**
719 Code for statistical analyses is provided as part of the supplement. Custom R scripts and their
720 input data for signature analysis are available on GitHub at https://github.com/HLee-
721 Six/colon_microbiopsies. All other code is available from the authors on request.
722
723 **Calling substitutions**
724 Substitution calling was broken down into three steps: mutation discovery; filtering to
725 produce a list of clean sites; and genotyping, where the presence or absence of every
726 mutation in every sample is evaluated.
727
728 First, mutations were initially discovered using the Cancer Variants through Expectation
729 Maximisation (CaVEMan) algorithm[46]. CaVEMan uses a naïve Bayesian classifier to derive
730 the probability of all possible genotypes at each nucleotide. CaVEMan copy number options
731 were set to major copy number 5 and minor copy number 2 for normal clones, as in our
732 experience this maximises sensitivity. The algorithm was run using an unmatched normal in
733 order to be able to derive phylogenies: had another sample from the same individual been
734 treated as a matched normal, early embryonic mutations would have been treated as germline
735 and discarded, resulting in incorrect trees.
736
737 Second, a number of post-processing filters were applied. These included filtering against a
738 panel of 75 unmatched normal samples to remove common single nucleotide polymorphisms,
739 post-processing as described previously[32] and two filters (only applied to whole genome
740 sequencing data) designed to remove mapping artefacts associated with BWA-MEM: the
741 median alignment score of reads supporting a mutation should be greater than or equal to
742 140, and fewer than half of these reads should be clipped. The library preparation protocol for
743 microbiopsies produced shorter library insert sizes than standard methods. Reads could
744 therefore overlap, resulting in double counting of mutant reads. Fragment-based statistics
745 were generated to prevent the calling of variant supported by a low number of fragments.
746 Variants were annotated by ANNOVAR[47] and fragment-based statistics (fragment coverage,
747 number of fragments supporting the variant, fragment-based allele fraction) were calculated
748 for each variant after the exclusion of marked PCR duplicates. In the rare event of
749 discordance in the called base at the variant position between overlapping paired-end reads,

750    the base with the highest quality score was selected. Fragment-based statistics were
751    calculated separately for high quality fragments (alignment score ≥ 40 and base scores ≥ 30).
752    Variants supported by at least three high quality fragments were retained and used for the
753    next stage of variant filtering. Inspection of variants specific to LCM experiments revealed
754    that the vast majority were present within inverted repeats capable of forming hairpin
755    structures, that they were supported by reads with very similar alignment start position (and
756    so not marked as PCR duplicates), and were primarily located close to the alignment start
757    within the supporting reads. Commonly these variants coincided with other proximal variants
758    (1-30 bp), but filtering based on variant proximity would also remove actual kataegis events.
759    *In silico* modelling of the potential hairpin showed that the variants were aligning to each
760    other in the stem of the structure, but could not form a base pair, while all other bases could.
761    The artefacts are likely the consequence of erroneous processing of cruciform DNA (existing
762    either prior to DNA isolation or formed during library preparation) by the enzymatic
763    digestion protocol applied. We have considered modelling the hairpin structures to filter these
764    variants, but given the fact that read clustering (i.e., similar alignment position) serves as a
765    hallmark for these artefacts, we opted to use the proximity of the variant to the alignment
766    start, and the standard deviation (SD) and median absolute deviation (MAD) of the variant
767    position within the supporting reads, as features for filtering. These statistics were calculated
768    separately for positive and negative strand aligned reads. In case the variant was supported by
769    a low number of reads (i.e., 0-1 reads) for one of the strands, the filtering was based only on
770    the statistics generated for the other strand. Per variant, if one of the strands had too few reads
771    supporting, it was required for the other strand that either: (I) there should be ≤ 90%
772    supporting reads to report the variant within the first 15% of the read starting from the
773    alignment start, or (II) the statistics MAD > 0 and SD > 4. Per variant, if both strands were
774    supported by sufficient reads it was required for both strands separately that either: (I) there
775    should be ≤ 90% supporting reads to report the variant within the first 15% of the read, (II)
776    the statistics MAD > 2 and a SD > 2, or (III) that the other strand should have the statistics
777    MAD > 1 and SD > 10 (i.e., the variant is retained if the other strand demonstrates strong
778    measures of variance). In our experience, the proposed strategy vastly reduces the number of
779    artefactual variants while retaining all other variants, as assessed by running the last filtering
780    step on WGS data from non-LCM experiments.
781
782    Third, mutations were genotyped in every sample. A pileup of all the samples from a given
783    individual was constructed, counting the number of mutant and wild type reads in every
784    sample over every site that had been called in any sample from that person. Only reads with a
785    mapping quality of 30 or above and bases with a base quality of 30 or above were counted.
786    After applying these filters, mutations were genotyped based on the number of mutant and
787    wild type reads at each locus. Mutations were called based on a variant allele fraction (VAF)
788    > 0.2, a depth > 7, and at least 4 mutant reads. If the depth over a locus was less than seven in
789    a given sample, or if there was more than one mutant read but the other criteria were not met,
790    the genotype was set to NA for tree construction purposes. Loci that were set to NA in more
791    than one third of the samples were removed for construction of the phylogeny. Positions were
792    called as germline if they were either called as present or NA in all of the samples from a
793    given individual.
794

| Mutations called against unmatched normal with CaVEMan algorithm | → | Filtered against 75 unmatched normals | → | Post processing filters as described in Nik-Zainal et al 2012 | → | Median alignment score of reads (ASMD) >=140 | → | Fewer than half the reads should be clipped (CLPM=0) | → | Fragment based and cruciform filters | → | Pileup: count mt and wt reads in all samples for subs called in any sample |

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

As for substitutions, calling of indels was broken down into mutation discovery, filtering, and genotyping. Mutations were called with the Pindel algorithm[48] using an unmatched normal. Post processing filters were applied as in Nik-Zainal et al.[32], and the number of mutant and wild-type reads was tabulated as above. The same dataset-specific filters were applied as for substitutions. Indels were then genotyped based on a VAF>0.2, a depth of at least 10, and support of at least 5 mutant reads.

**Calling structural variants**

Genomic rearrangements were called using the BRASS algorithm[41] (https://github.com/cancerit/BRASS). Abnormally paired read pairs from WGS were grouped and filtered by read remapping. Read pair clusters with ≥50% of the reads mapping to microbial sequences were removed, as were rearrangements where the breakpoint could not be reassembled. Candidate breakpoints were matched to copy number breakpoints defined by ASCAT (see below) within 10kb. Only structural variants where the two breakpoints were more than 1000 base pairs apart were considered. Structural variants were called against a matched normal skin or blood sample when available and against another crypt from the same individual with good coverage when not.

**Calling copy number**

Copy number changes were called using the Allele-Specific Copy number Analysis of Tumours (ASCAT) algorithm[49]. The same matched normal sample was used as for calling structural variants. For additional validation of copy number changes in normal colon, the QDNAseq algorithm[50] was run. ASCAT uses both the read depth and ratios of heterozygous single nucleotide polymorphisms to determine an allele-specific copy number, while the QDNAseq relies solely on variations in sequencing coverage. To call amplifications and deletions in the colonic microbiopsy cohort, only those that were both called by ASCAT and showed a clear departure from the background log2ratio by QDNAseq were retained. To call copy neutral loss of heterozygosity in this cohort, all such events called by ASCAT were checked visually on Jbrowse[51] to verify an imbalance of parental snps. Only crypts with >10X coverage, for which copy number changes could be reliably detected, were used.

**Detection of driver variants and positive selection**

Driver mutations were detected both through an unbiased dNdS method and through manual annotation. For these analyses, the CaVEMan and Pindel calls were used without post-processing filters in order to maximise our sensitivity. All putative driver variants were visually inspected using Jbrowse[51], and so we could afford a higher false positive rate in the mutation discovery phase.

dNdScv[52] was used to conduct three tests: first, using only the whole genome sequencing data, an analysis of selection over all genes; second, using combined whole genome and targeted sequencing data, over all the genes covered by the bait-set; and finally, using again

839   this combined dataset, over 90 selected cancer genes (appendix). R code for this analysis is
840   included in the supplementary information.
841
842   Manual annotation of driver variants based on prior knowledge complemented this. A list of
843   90 colorectal cancer genes (appendix) curated from the literature that were also covered by
844   the bait-set were intersected with the list of substitutions and indels from combined whole
845   genome and targeted sequencing. Mutations were annotated as putative drivers if they were
846   either missense mutations that fell in an oncogene hotspot (based on visualisation of the
847   distribution of mutations in the gene on COSMIC[53]), or if they were truncating mutations that
848   fell in a tumour suppressor gene.
849
850   Structural variants that might act as drivers were assessed by intersection of genes involved
851   in each structural variant with the twelve genes involved in gene fusions that have been
852   reported in colorectal cancer in COSMIC (*VTI1A, TCF7L2, TPM3, NTRK1, PTPRK, RSPO3,*
853   *ETV6, NTRK3, EIF3E, RSPO2, C2orf44*, and *ALK*). No fusion genes were found. None of the
854   genes involved in structural variants in our data overlapped with the list of 90 cancer genes
855   used for assessing substitutions and indels, and nor were there any genes that were affected
856   by more than one structural variant. No high-level copy number amplifications were observed
857   and there were no homozygous deletions.
858
859   **Estimation of frequency of driver mutations in cancer**
860   Publically-available colorectal cancer mutation calls were obtained from The Cancer Atlas
861   Network[43]. Driver mutations were annotated manually in the same way as in our dataset: only
862   mutations that fell in the 90 genes that we had selected were considered, and they were
863   annotated as putative drivers if they were either missense mutations that fell in an oncogene
864   hotspot (based on visualisation of the distribution of mutations in the gene on COSMIC[53]), or
865   if they were truncating mutations that fell in a tumour suppressor gene.
866
867   **Construction of phylogenies**
868   Phylogenies are used in this analysis for timing mutations. The most informative branches in
869   this case are the long branches shared by a small number of crypts, which are very robust to
870   all tree construction methods. Trees were built using maximum parsimony using substitutions
871   called as described above. For every individual, the input matrix of mutation calls was
872   bootstrapped 100 times. Phylogenies were constructed for each replicate using the Wagner
873   method of the Mix programme from the Phylip suite of tools[54]. The consensus of all the
874   phylogenies constructed was used.
875
876   The phylogenies were validated using the indel calls. To do this, the same procedure as for
877   substitutions was followed for indel matrices. As there were fewer indels than substitutions,
878   nodes in indel phylogenies were generally reconstructed with lower confidence than in
879   substitution phylogenies, but they broadly agree. 85% of nodes reconstructed with >=90%
880   confidence in the indel tree were present with exactly the same set of descendants in the
881   substitution trees.
882
883   The phylogeny inference programme used provided the topology of the tree but not the
884   assignment of mutations. Mutations from the input matrix of genotypes therefore have to be
885   re-assigned to branches. In order to assign a set of mutation calls with no false negative and
886   no false positives to a tree, each branch of the tree was considered in turn. If a mutation was
887   called in all the descendants of a given branch, and in no samples that were not descendants
888   of the branch, mutations were assigned to that branch.

889
890   Some colonic microbiopsies suffered from low coverage and stromal contamination. For this
891 reason, we did not expect mutations to fit the tree perfectly, as a mutation that was truly
892 present in a colony might be missed if too few supporting reads are found. Mutations were
893 only assigned to the tree in order to determine the mutational processes active at a particular
894 time. We reasoned that it was preferable to assign only mutations that fit the tree perfectly
895 and adjust the branch lengths based on the power to call mutations at a given branch, rather
896 than attempting to assign mutations that fit the tree imperfectly. Using the clonality and
897 coverage of all descendants of a branch, the proportion of true substitutions or indels on the
898 branch that would be first discovered (whether by CaVEMan or Pindel) and then genotyped
899 as present according to the criteria described above was calculated. The observed branch
900 length was then adjusted by dividing by this proportion. This was done for both substitutions
901 and indels, but not for structural variants and for larger copy number changes due to a lack of
902 data: most branches have no large variants and so could not be extended appropriately.
903 Rearrangements and copy number changes were assigned to phylogenies manually.
904
905 **Extraction of mutational signatures**
906 Mutational signatures were extracted using the mutations assigned to every branch of a
907 phylogeny as a 'sample'. This allows better discrimination of mutational processes that may
908 occur at different times within the same cell. Mutations were categorised following the
909 method used by the Mutational Signatures working group of the Pan Cancer Analysis of
910 Whole Genomes (PCAWG)[1]. Single base substitutions were categorised into 96 classes
911 according the identity of the pyrimidine mutated base pair, and the base 5' and 3' to it.
912 Doublet base substitutions were categorised into 78 classes according to the identity of the
913 reference and alternative bases. Indels were classified according to whether they were an
914 insertion or a deletion, the identity of the inserted/deleted base, the length of the
915 mononucleotide tract in which they occurred, or the degree of homology with the
916 surrounding sequence into 83 classes (Fig. 1a).
917
918 Signatures were extracted using a hierarchical Dirichlet Process[55,56]. Code and the input
919 mutations are provided at https://github.com/HLee-Six/colon_microbiopsies. First, the
920 algorithm was conditioned on the set of mutational signatures that have found to be operative
921 in colorectal cancers in PCAWG[1]: SBS1, SBS2, SBS3, SBS5, SBS13, SBS16, SBS17a,
922 SBS17b, SBS18, SBS25 (included although it is not found in colorectal cancer because the
923 similarity with the mutational profile with crypts from one individual had been previously
924 noted), SBS28, SBS30, SBS37, SBS40, SBS41, SBS43, SBS45, SBS49, DBS, DBS3, DBS4,
925 DBS6, DBS7, DBS8, DBS9, DBS10, DBS11, ID1, ID2, ID3, ID4, ID5, ID6, ID7, ID8, ID10,
926 and ID14. This allows simultaneous discovery of new signatures and matching to known
927 ones. Nine single base substitution (SBS), two doublet base substitution (DBS), and five
928 indel (ID) signatures were discovered (Extended Data Fig. 2). Despite pre-conditioning,
929 signatures that were perfectly correlated in all samples were still amalgamated. This
930 occurred, for example, with signatures 1, 5, and 18. Therefore, expectation maximisation was
931 used to deconvolute all HDP signatures into known PCAWG signatures. If a signature
932 reconstituted from the components that expectation maximisation extracted (only including
933 PCAWG signatures that accounted for at least 10% of mutations in each sample to avoid
934 over-fitting) had a cosine similarity to the HDP signature of more than 0.95, the signature
935 was presented as its expectation maximisation deconvolution. Three HDP signatures met
936 these criteria: the HDP SBS1 signature was deconvoluted into a mixture of PCAWG SBS1,
937 PCAWG SBS5, and PCAWG SBS18; the HDP DBSA was deconvoluted in PCAWG DBS2,
938 PCAWG DBS4, PCAWG DBS6, PCAWG DBS9, and PCAWG DBS11; and the HDP IDC

939    was deconvoluted into PCAWG ID1, PCAWG ID2, and PCAWG ID5 (Extended Data Fig.
940    3). To test the robustness of this signature analysis, other signature extraction methods were
941    used: HDP with no pre-conditioning, the non-negative matrix factorisation (NNMF) method
942    used by Blokzijl and colleagues[4], and a version of the NNMF algorithm used by Alexandrov
943    and colleagues[1]. These all produced comparable results (Extended Data fig. 4).
944
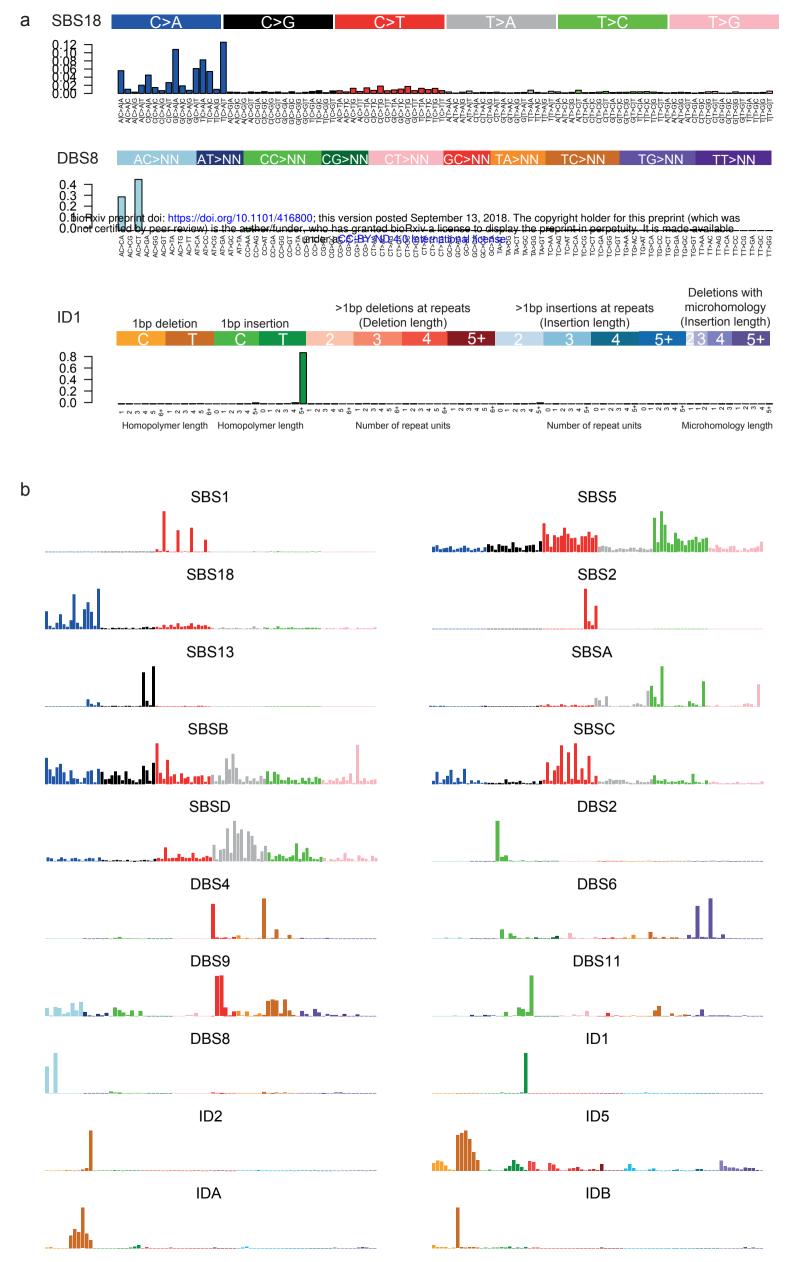
945    **Statistical analyses**
946    All statistical analyses were performed in R (Supplementary Results).
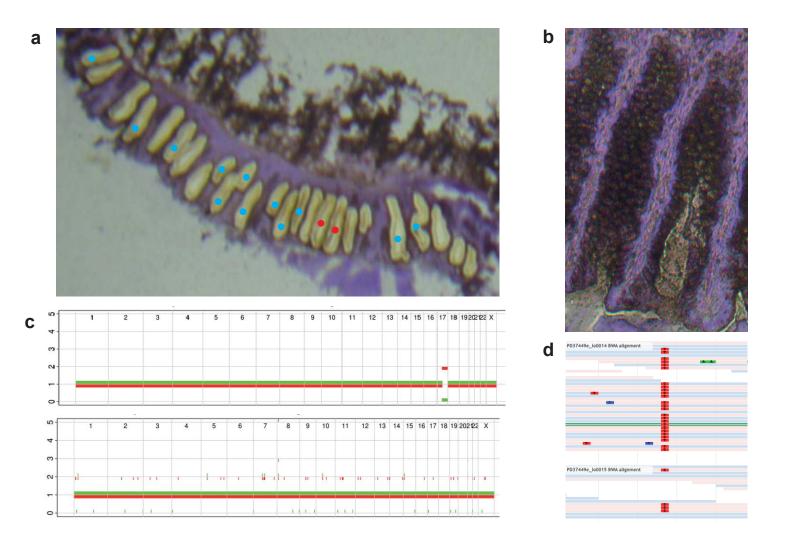
947
948
949    **REFERENCES FOR SUPPLEMENTARY METHODS**
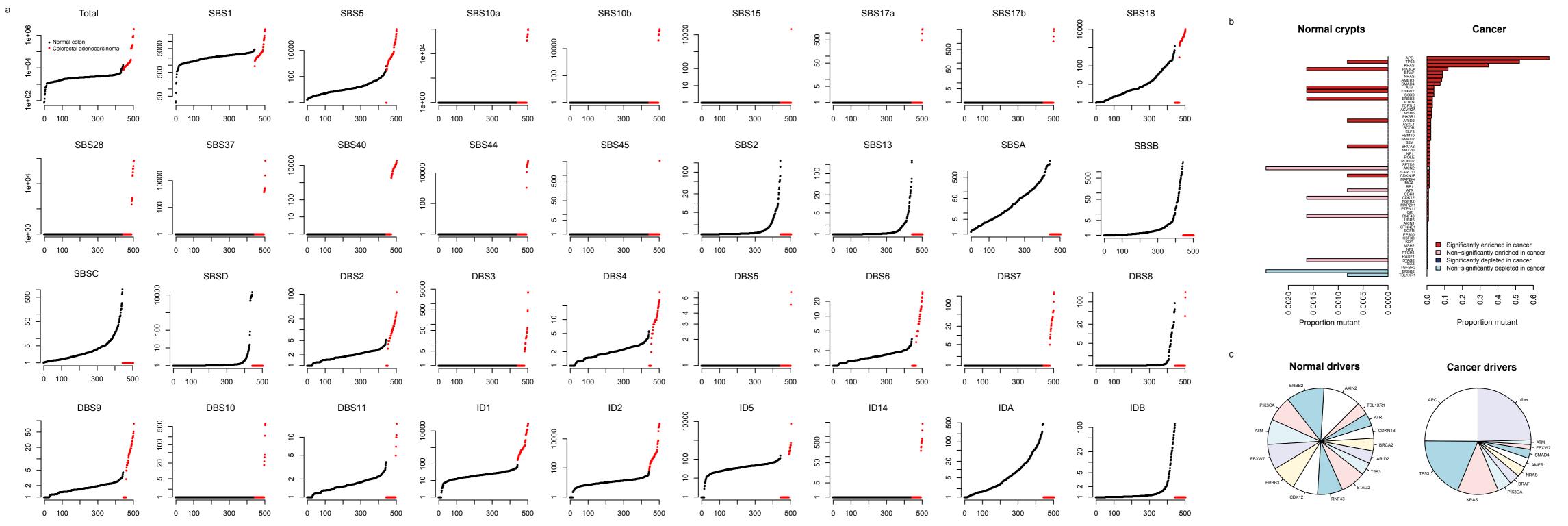950

951    46. Jones, D. et al. cgpCaVEManWrapper: Simple execution of CaVEMan in order to
952        detect somatic single nucleotide variants in NGS data. *Curr. Protoc. Bioinformatics*
953        **56**, 15.10.1–15.10.18 (2016).
954    47. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic
955        variants from high-throughput sequencing data. *Nucleic acids research* **38**, e164,
956        (2010).
957    48. Raine, K. M. et al. cgpPindel: identifying somatically acquired insertion and deletion
958        events from paired end sequencing. *Curr. Protoc. Bioinformatics* **52**, 15.7.1–15.7.12
959        (2015).
960    *49.* Van Loo, P. et al. Allele-specific copy number analysis of tumors. *Proc. Natl Acad.*
961        *Sci. USA* **107**, 16910–16915 (2010).
962    *50.* Scheinin, I. et al. DNA copy number analysis of fresh and formalin-fixed specimens
963        by shallow whole-genome sequencing with identification and exclusion of
964        problematic regions in the genome assembly. *Genome Research*, **24**, 2022–2032
965        (2014).
966    51. Buels R *et al*. JBrowse: a dynamic web platform for genome visualization and
967        analysis. *Genome Biology* doi: 10.1186/s13059-016-0924-1 (2016).
968    52. Martincorena, I. et al. Universal patterns of selection in cancer and somatic tissues.
969        *Cell* **171**, 1029–1041 (2017).
970    53. Forbes S.A. et al. COSMIC: somatic cancer genetic sat high-resolution. *Nucleic Acids*
971        *Res.* **45**, D777-D783 (2017).
972    54. Felsenstein, J. PHYLIP — Phylogeny Inference Package (Version 3.2). *Cladistics* **5**,
973        164–166 (1989).
974    55. Roberts, N. *Patterns of somatic genome rearrangement in human cancer.* PhD thesis,
975        Univ Cambridge, UK (Wellcome Trust Sanger Institute, 2018).
976    56. Nicola Roberts, R pkg for Hierarchical Dirichlet Process,
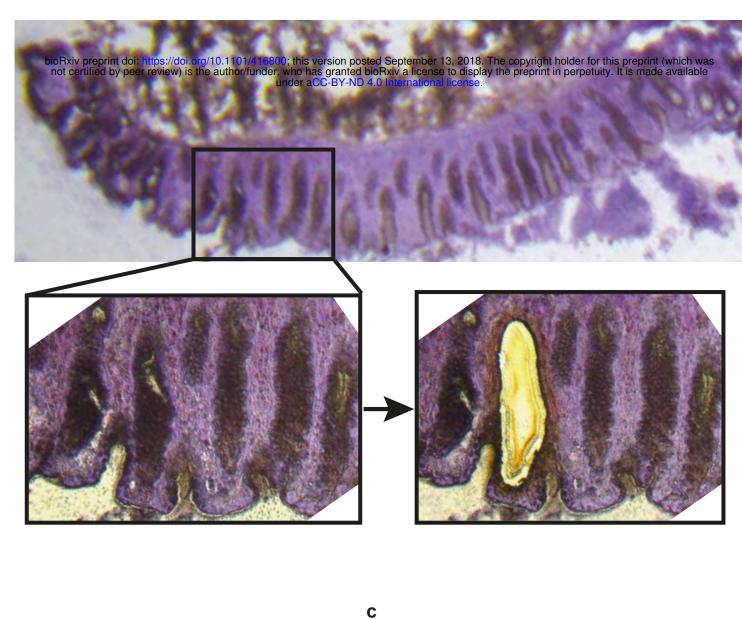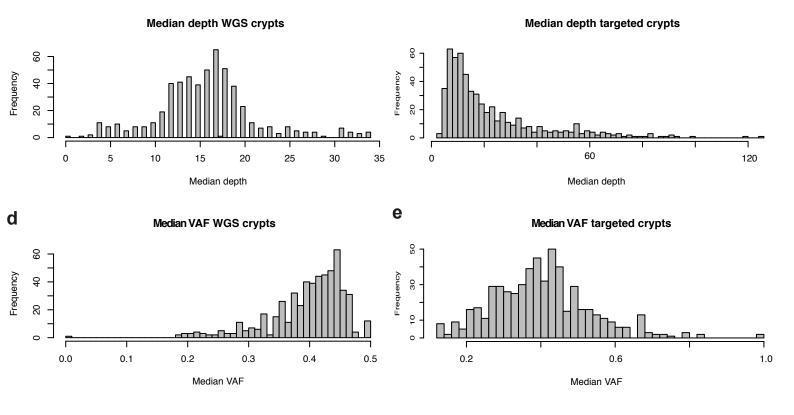977        https://github.com/nicolaroberts/hdp (Accessed August 2018).

**a**

SBS18

DBS8

ID1

**b**

SBS1 SBS5

SBS18 SBS2

SBS13 SBSA

SBSB SBSC

SBSD DBS2

DBS4 DBS6

DBS9 DBS11

DBS8 ID1

ID2 ID5

IDA IDB

**a**

hdp sig1 · hdp sig2 · hdp sig3
hdp sig4 · hdp sig5 · hdp sig6
hdp sig7 · hdp sig8 · hdp sig9

**b**

nmf sig1 · nmf sig2 · nmf sig3
nmf sig4 · nmf sig5 · nmf sig6
nmf sig7 · nmf sig8 · nmf sig9

**c**

mutpat sig1 · mutpat sig2 · mutpat sig3
mutpat sig4 · mutpat sig5 · mutpat sig6
mutpat sig7 · mutpat sig8 · mutpat sig9

## SBS1

Ileum: 12.7 (CI95 10.5–14.7)
Right: 12.9 (CI95 11.4–14.3)
Transverse: 16 (CI95 14.7–17.6)
Left: 12.7 (CI95 11.3–14.4)

## SBS5

Ileum: 16.4 (CI95 13.7–18.9)
Right: 16.2 (CI95 14.3–18.2)
Transverse: 20.4 (CI95 18.3–22.2)
Left: 16.4 (CI95 14.3–18.3)

## SBS18

Ileum: 5.13 (CI95 3.77–6.56)
Right: 7.48 (CI95 6.37–8.54)
Transverse: 6.99 (CI95 5.9–8.03)
Left: 5.38 (CI95 4.29–6.49)

## ID1

Ileum: 0.218 (CI95 0.105–0.344)
Right: 0.378 (CI95 0.293–0.467)
Transverse: 0.372 (CI95 0.287–0.46)
Left: 0.31 (CI95 0.223–0.405)

## ID2

Ileum: 0.0885 (CI95 0.0123–0.162)
Right: 0.171 (CI95 0.119–0.222)
Transverse: 0.18 (CI95 0.127–0.231)
Left: 0.124 (CI95 0.0732–0.18)

## ID5

Ileum: 0.402 (CI95 0.203–0.615)
Right: 0.675 (CI95 0.523–0.834)
Transverse: 0.68 (CI95 0.525–0.84)
Left: 0.547 (CI95 0.396–0.704)

*AXIN2*
loss

*PIK3CA
E542K*

Insufficient coverage

Insufficient coverage

Insufficient coverage

| | |
|---|---|
| sbs_X0 | |
| SBS1 | DBS4 |
| SBS5 | DBS6 |
| SBS18 | DBS9 |
| SBS2 | DBS11 |
| SBS13 | DBS8 |
| SBSD | id_X0 |
| SBSA | ID1 |
| SBSB | ID2 |
| SBSC | ID5 |
| dbs_X0 | IDA |
| DBS2 | IDB |

a

b



**Timing of copy number changes**