

Somatic mutation landscapes at single-molecule resolution

Federico Abascal¹, Luke M. R. Harvey^{1,#}, Emily Mitchell^{1,2,#}, Andrew R. J. Lawson^{1,#},
Stefanie V. Lensing^{1,#}, Peter Ellis^{1,3,#}, Andrew J. C. Russell¹, Raul E. Alcantara¹, Adrian
Baez-Ortega¹, Yichen Wang¹, Eugene Jing Kwa¹, Henry Lee-Six¹, Alex Cagan¹, Tim H. H.
Coorens¹, Michael Spencer Chapman¹, Sigurgeir Olafsson¹, Steven Leonard¹, David Jones¹,
Heather E. Machado¹, Megan Davies², Nina F. Øbro^{2,4}, Krishnaa T. Mahubani^{4,5,6}, Kieren
Allinson⁷, Moritz Gerstung⁸, Kourosh Saeb-Parsy^{5,6}, David G. Kent^{2,9}, Elisa Laurenti^{2,4},
Michael R. Stratton¹, Raheleh Rahbari¹, Peter J. Campbell^{1,4}, Robert J. Osborne^{1,10,*}, Iñigo
Martincorena^{1,*}.

These authors contributed equally

* Corresponding authors: *r.osborne@biofidelity.com* (R.J.O.), *im3@sanger.ac.uk* (I.M.)

Affiliations:

¹ Wellcome Sanger Institute, Hinxton CB10 1SA, UK.

² Wellcome - MRC Cambridge Stem Cell Institute, Cambridge Biomedical Campus,
Cambridge CB2 0AW, UK.

³ Current address: Inivata, Glenn Berge Building, Babraham Research Campus, Babraham,
Cambridge, CB22 3FH, UK

⁴ Department of Haematology, University of Cambridge, Cambridge CB2 2XY, UK.

⁵ Department of Surgery, University of Cambridge, Cambridge CB2 0QQ, UK.

⁶ NIHR Cambridge Biomedical Research Centre, Cambridge Biomedical Campus,
Cambridge CB2 0QQ, UK.

⁷ Cambridge Brain Bank, Division of the Human Research Tissue Bank, Box 235, Level 5,
Addenbrooke's Hospital, Hills Rd, Cambridge, CB2 0QQ, UK.

⁸ European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI),
Hinxton CB10 1SD, UK.

⁹ York Biomedical Research Institute, Department of Biology, University of York, York
YO10 5DD, UK.

¹⁰ Current address: Biofidelity, 330 Cambridge Science Park, Milton Road, Cambridge, CB4
0WN, UK

Somatic mutations drive cancer development and may contribute to ageing and other diseases^{1,2}. Yet, the difficulty of detecting mutations present only in single cells or small clones has limited our knowledge of somatic mutagenesis to a minority of tissues. To overcome these limitations, we introduce nanorate sequencing (NanoSeq), a new duplex sequencing protocol with error rates <5 errors per billion base pairs in single DNA molecules from cell populations. This rate is two orders of magnitude lower than typical somatic mutation loads, enabling the study of somatic mutations in any tissue independently of clonality. We exploited this single-molecule sensitivity to study somatic mutations in non-dividing cells across several tissues, comparing stem cells to differentiated cells and studying mutagenesis in the absence of cell division. Differentiated cells in blood and colon displayed remarkably similar mutation loads and signatures to their corresponding stem cells, despite mature blood cells having undergone considerably more divisions. We then characterised the mutational landscape of post-mitotic neurons and polyclonal smooth muscle, confirming that neurons accumulate somatic mutations at a constant rate throughout life without cell division, with similar rates to mitotically-active tissues. Altogether, our results suggest

51 **that mutational processes independent of cell division are important contributors to**
52 **somatic mutagenesis. We anticipate that the ability to reliably detect mutations in single**
53 **DNA molecules could transform our understanding of somatic mutagenesis and enable**
54 **non-invasive studies on large-scale cohorts.**

55 56 **Introduction**

57
58 Somatic mutations occur in our cells as we age. Because most somatic mutations are present
59 in small groups of cells or even in single cells, studying somatic mutagenesis has been
60 challenging, requiring special approaches. This includes ultra-deep sequencing of small
61 biopsies³⁻⁵, laser microdissection⁶⁻⁸, isolation of single-cells followed by in vitro expansion
62 into organoids or colonies⁹⁻¹¹, and single-cell sequencing¹²⁻¹⁴. While these technologies are
63 changing our understanding of somatic mutagenesis, the error rates of single-cell approaches
64 have, until recently¹⁵, been too high¹⁶, and other approaches are limited to mitotically-active
65 cell types.

66
67 As a result of these technical limitations, the rates and patterns of somatic mutation across
68 most human cell types remain underexplored. This is especially the case for non-dividing
69 cells, including differentiated cells that make the bulk of mitotically-active tissues and are
70 responsible for tissue function, and post-mitotic tissues, such as cortical neurons or cardiac
71 muscle, which are of particular interest in human ageing, neurodegeneration and
72 cardiovascular disease. Post-mitotic tissues can also inform on the contribution of cell
73 division and DNA replication to somatic mutation in human tissues. To address these
74 questions, here we present a new sequencing protocol that enables the study of somatic
75 mutations in any tissue or cell population by reliably detecting somatic mutations in single
76 DNA molecules.

77 78 **Nanorate sequencing**

79
80 Several protocols have been developed to increase the accuracy of standard sequencing by
81 barcoding individual molecules of DNA and sequencing each molecule multiple times,
82 reducing error rates by single-molecule consensus¹⁷. The most accurate approaches use
83 duplex consensus sequencing^{18,19}, sequencing copies of both strands of a DNA molecule to
84 remove sequencing errors (present in individual reads) and PCR errors (present in copies of
85 one of the two strands) (**Fig. 1a**). Duplex sequencing has a theoretical error rate $<10^{-9}$
86 errors/bp, the probability of two early and complementary PCR errors in both strands¹⁷.
87 Given that this rate is lower than the typical mutational load of human tissues, it raises the
88 possibility of quantifying somatic mutation rates in genetically-heterogeneous samples, by
89 detecting somatic mutations on single DNA molecules. This is the rationale of BotSeqS, a
90 whole-genome duplex sequencing protocol²⁰ (**Fig. 1a**). In practice, however, mapping errors
91 and some library preparation artefacts can violate the assumed independence of both
92 strands^{20,21}. The actual error rates of duplex sequencing protocols have remained difficult to
93 measure due to the lack of control samples with low and known mutation rates¹⁷.

94
95 To evaluate the performance of BotSeqS, we used samples of cord blood, comparing
96 BotSeqS on bulk granulocytes from a neonate to standard sequencing of 100 single-cell-
97 derived colonies from two neonates as a control. On average, single-cell-derived colonies had
98 66 mutations per cell, dominated by C>T mutations at CpG sites. In contrast, BotSeqS
99 estimated 1,240 mutations per diploid genome, dominated by C>A and C>G (**Fig. 1b, c**).
100 Analysing the distribution of substitutions across BotSeqS reads revealed a large excess of

101 G>T/C and C>T substitutions near the 5' ends of DNA fragments, and an imbalance over the
102 complementary C>A/G and G>A substitutions affecting the entire read length (**Fig. 1d**,
103 **Extended Data Figs. 1** and **2**). These imbalances are incompatible with real mutations and
104 reflect errors introduced during library preparation²² (**Methods, Supplementary Note 1**). We
105 found the same imbalances, with a much larger C>T component, in the original BotSeqS
106 publication²⁰ (**Fig. 1d**). Extensive trimming of read ends only partially alleviated these errors
107 (**Extended Data Fig. 2**). Overall, we estimate that BotSeqS introduced approximately 1,200
108 errors per diploid genome in our samples (i.e. $\sim 2 \times 10^{-7}$ errors/bp).
109

110 Based on the error patterns, we reasoned that end repair was likely responsible for most
111 errors, by converting DNA damage in single-strands of DNA into double-stranded errors
112 (**Fig. 1e** and **Extended Data Fig. 1c, d**). To solve this, we developed NanoSeq, a protocol
113 that prevents copying errors between strands by avoiding end repair and by blocking nick
114 extension. First, we replaced sonication and end repair with restriction enzyme fragmentation
115 (**Fig. 1e, Methods, Supplementary Table 3, Supplementary Note 2**). Although restriction
116 enzymes provide partial coverage of the genome (29% using HpyCH4V), the fraction
117 covered is sufficiently random to accurately estimate mutation rates and signatures. They also
118 enable the generation of NanoSeq libraries from as little as 1 ng of DNA (**Methods**).
119 Alternatively, we show that sonication followed by exonuclease blunting can be used for
120 applications requiring whole-genome coverage (**Methods, Supplementary Note 3**,
121 **Extended Data Fig. 3**). Second, we introduced non-A dideoxynucleotides (ddBTPs) during
122 A-tailing, to avoid errors from nick extension (**Fig. 1e, Methods, Extended Data Fig. 1e**,
123 **Supplementary Note 4**). Adapters with sufficiently diverse random barcodes were used to
124 create single-molecule-derived read families (**Supplementary Note 5**).
125

126 If duplicate rates are not optimised, duplex sequencing approaches can suffer from low
127 efficiency due to suboptimal read family sizes²⁰. We use mathematical modelling of family
128 sizes and qPCR quantification of the library to maximise the duplex coverage independently
129 of the amount of input DNA (**Methods, Extended Data Fig. 4a-d**). A robust bioinformatic
130 pipeline was also developed to avoid false positive mutation calls from mapping errors and
131 from low-level DNA contamination (**Extended Data Fig. 4e, f, Methods, Supplementary**
132 **Note 6**), and to distinguish germline from somatic mutations.
133

134 Applying NanoSeq to cord blood granulocytes yielded an estimated mutation rate of 109
135 mutations per cell (95% Poisson confidence intervals 95-125; **Fig. 1g**). The small difference
136 with the colonies could be due to NanoSeq errors, a higher mutation burden in granulocytes
137 than in cord blood stem cells, or both. Consistent with most mutations detected by NanoSeq
138 being genuine, no substitution imbalances were detected in the NanoSeq calls (**Fig. 1d**) and
139 no significant differences were found between the mutational spectra of colonies and
140 granulocytes (**Fig. 1c, Methods**). As an additional low-burden control, we applied NanoSeq
141 to a sperm sample from a 21-year-old donor. Seven NanoSeq replicates of the sperm sample
142 yielded low mutation burdens, with ~ 52 mutations per haploid sperm cell (1.8×10^{-8}
143 mutations/bp or ~ 2.5 mutations/year/cell), consistent with current estimates of the mutation
144 rate in the paternal germline from trio studies^{23,24} (**Fig. 1f**). Together, the sperm and cord
145 blood data indicate that the error rate of NanoSeq is lower than 5×10^{-9} errors/bp (<30 errors
146 per diploid genome), two orders of magnitude lower than the BotSeqS error rate and the
147 somatic mutation load of most human tissues studied to date. Analysis of insertions and
148 deletions (indels) also revealed an indel error rate $< 3 \times 10^{-9}$ errors/bp (**Methods, Extended**
149 **Data Fig. 5c, Supplementary Note 8**).
150

151 The extremely low error rate of NanoSeq, in the nano range, enables the reliable detection of
152 somatic mutations in single DNA molecules, opening the door to the study of somatic
153 mutations in any tissue or cell population. We take advantage of this unprecedented ability to
154 study non-dividing cells across four tissues, addressing two elusive questions in the field of
155 somatic mutagenesis: the difference in mutation rates between stem cells and terminally-
156 differentiated cells in mitotically-active tissues, and the rates and patterns of mutation in post-
157 mitotic tissues.

158

159 **Mutation burden in stem and differentiated cells**

160

161 Due to technical limitations, most of our knowledge of somatic mutagenesis is restricted to
162 stem or proliferating cells. Since stem cells are believed to be better protected against
163 mutations²⁵, differentiated cells could conceivably have higher mutational loads and
164 undescribed mutational signatures¹⁴.

165

166 We first addressed this question in the haematopoietic system, comparing mature
167 granulocytes to haematopoietic stem and multipotent progenitor cells (HSC/MPPs)
168 (**Methods**). The haematopoietic system is organised hierarchically, with a heterogeneous
169 pool of slow-cycling stem cells sustaining the production of large numbers of differentiated
170 cells through the extensive proliferation of intermediate progenitor cells (**Fig. 2a**). HSCs are
171 estimated to divide around once a year and conservative estimates suggest that an average of
172 over 28 cell divisions must separate stem cells from differentiated cells to explain the
173 production of $\sim 10^{14}$ mature cells per year (**Fig. 2a, Supplementary Note 9**). As a result, a
174 considerably higher mutation burden and mutational signatures associated with the
175 proliferation of progenitors may be expected in granulocytes.

176

177 We used NanoSeq to sequence 18 samples of granulocytes from 9 healthy donors, ranging
178 from 20 to 80 years of age (**Supplementary Table 1, 2**). We compared these data to standard
179 whole-genome sequencing of 60 single-cell derived HSC/MPPs colonies from 6 donors
180 (**Extended Data Fig. 6a, Supplementary Table 1, 2**) and published data from 110 colonies
181 from one donor²⁶ (**Methods**). These data revealed remarkably similar mutation burdens in
182 terminally-differentiated granulocytes and HSC/MPPs (**Fig. 2b**). Linear mixed-effect
183 regression yielded indistinguishable slopes for HSC/MPPs colonies and granulocytes
184 ($P=0.92$), with a joint estimate of ~ 19.9 mutations/year (CI95% 18.3-21.4, **Methods**,
185 **Supplementary Table 8**). The excess of mutations in granulocytes over HSC/MPPs was
186 estimated to be ~ 51 mutations and not significantly different from zero (CI95%: -14-120,
187 $P=0.13$, **Methods, Supplementary Table 8**). Their mutational spectra were also largely
188 similar (cosine similarity 0.98, **Fig. 2c**).

189

190 The observation that a considerable increase in cell divisions does not cause a proportional
191 increase in mutation burden suggests that replication errors cannot be responsible for more
192 than a small minority of mutations in HSC/MPPs (**Supplementary Note 9**). A caveat for this
193 comparison is that HSC/MPP colonies successfully grown in vitro may not reflect the
194 mutation rate of the more quiescent HSCs responsible for long-term maintenance of the
195 haematopoietic system. However, a similar conclusion can be drawn from the granulocyte
196 data alone. The strong linear relationship with age and the small intercept for granulocytes
197 alone (142.1 mutations, CI95%: -115.3-414.2, compared to the slope of ~ 19.8
198 mutations/year) suggests that the majority of the mutations observed in adult granulocytes
199 accumulated in the stem cells responsible for long-term maintenance, and that only a small

200 minority of mutations are accrued during transient proliferation and terminal differentiation
201 (**Supplementary Note 9**).

202

203 To extend the comparison of stem cells and differentiated cells to another tissue with a well-
204 understood stem cell organisation, we studied colonic epithelium. Estimates of the somatic
205 mutation rate in colonic stem cells are available from whole-genome sequencing of clonal
206 organoids derived from single Lgr5+ cells¹⁰ and from sequencing single laser-microdissected
207 colonic crypts⁶, which over time become clonally derived from a single stem cell²⁷. For three
208 previously-studied donors we compared standard whole-genome sequencing of
209 microdissected colonic crypts⁶ to NanoSeq data from single crypts or groups of crypts
210 (**Extended Data Fig. 6b, c**). This revealed similar estimates of mutation burden, despite the
211 time lag to clonality in standard sequencing of colonic crypts (**Fig. 2d**). Mutation burden and
212 signatures from differentiated cells in colonic epithelium were consistent with those found by
213 previous studies on colonic stem cells, with a dominance of SBS1, SBS5 and, in some
214 donors, a colibactin signature²⁸ (**Fig. 2e** and **Extended Data Fig. 6d**).

215

216 Overall, NanoSeq data on granulocytes and colonic epithelium yielded similar mutation
217 burdens and signatures to their corresponding stem cells. While larger studies will be needed
218 to identify subtler differences and to address this question in other cell types, these results
219 provide an early view into the somatic mutation landscape of two differentiated cell types.

220

221 **Mutagenesis in neurons and smooth muscle**

222

223 Cortical neurons are a prime example of a post-mitotic tissue. This makes them a key cell
224 type to study somatic mutagenesis in the absence of cell division, but also inaccessible to
225 traditional sequencing methods. Despite the technical challenges impeding progress, somatic
226 mutations in neurodegeneration have attracted considerable interest^{1,12,13,29}.

227

228 We applied NanoSeq to frontal cortex neurons from 8 healthy donors and 9 Alzheimer's
229 disease (AD) patients (**Supplementary Table 1**), using nuclei sorting with the *NeuN*
230 neuronal marker (**Methods, Extended Data Fig. 7a**). These data revealed a linear
231 accumulation of 17.1 substitutions (linear regression, CI95%:13.7-20.5) and 2.5 indels
232 (CI95%:1.7-3.3) per year, approximately constant throughout life (**Fig. 3a, b**,
233 **Supplementary Table 8**). This confirms that mutations accumulate in a clock-like fashion in
234 cortical neurons, in the absence of cell division, consistent with observations from single-cell
235 sequencing¹³.

236

237 A previous study using SNP-phased error-corrected single-cell sequencing reported three
238 signatures in neurons, one that increased linearly with age and two that did not¹³. The
239 spectrum found by NanoSeq and the mutation rate per year closely resemble the age-
240 associated signature in that study (cosine similarity 0.96; **Fig. 3a, c** and **Extended Data Fig.**
241 **7b, c**). The two other signatures, responsible for around 72% of all mutations reported in the
242 study (**Extended Data Fig. 7d**), appear exclusively in single-cell data and likely derive from
243 amplification errors or transient DNA damage. Consistent with this possibility, the dominant
244 signature in single-neuron data closely resembles a single-cell-specific signature reported in
245 vitro¹⁶ (cosine similarity 0.97, **Extended Data Fig. 7b**).

246

247 To better understand the mutational processes active in post-mitotic neurons, we performed
248 signature decomposition on NanoSeq data from neurons, granulocytes, colonic crypts and
249 smooth muscle (described below). Three signatures were extracted (**Fig. 3e**): signatures A

250 and C imperfectly resembled SBS5 (cosine similarity 0.80) and SBS16 (0.78), respectively,
251 while signature B closely matched SBS1 (C>T changes at CpG dinucleotides, cosine
252 similarity 0.96). It is conceivable that SBS5, which appears to be a ubiquitous signature in
253 normal tissues and cancer genomes³⁰, reflects a collection of co-occurring processes, rather
254 than a single mutational process, leading to some differences across tissues. The observation
255 in post-mitotic neurons of signatures resembling SBS5 and SBS16 suggests that these
256 common processes, whose aetiologies remain poorly understood, can occur independently of
257 cell division.

258

259 The mutational spectra from neurons (**Fig. 3c, d**) showed several interesting features. T>C
260 substitutions at ApT sites appear enriched in neurons and show strong transcriptional strand
261 biases (**Extended Data Fig. 8b, c**). Signature B (SBS1), which is believed to be caused by 5-
262 methylcytosine deamination and fixed during DNA replication, accumulates at a low rate
263 with age in neurons (2.5 substitutions per year, linear regression CI95% 0.9-4.1, $P = 0.005$;
264 **Extended Data Figs. 7e and 9a, b**). This suggests that 5-methylcytosine deamination can be
265 fixed in both DNA strands without cell division, possibly by DNA repair. Neurons also have
266 a higher proportion of indels than other tissues, with an unusual enrichment of indels longer
267 than 1bp in highly-expressed genes, a pattern that resembles a mutational process recently
268 described in cancer genomes³¹ (**Fig. 3d, f** and **Extended Data Fig. 9c, d**). In contrast to other
269 somatic tissues, neurons did not exhibit a clear association between expression levels and
270 substitution rates across genes (**Fig. 3g**) and the enrichment of mutations in heterochromatin
271 was weaker (**Fig. 3h** and **Extended Data Fig. 8a**).

272

273 Although the difference is small, AD donors showed a slightly lower substitution rate than
274 healthy donors (linear regression, 17.6 (CI95%:15.0-20.2) vs 19.9 (CI95% 16.8-23.0)
275 substitutions/year, $P = 0.0029$) (**Fig. 3i, Extended Data Fig. 7e, Supplementary Table 8**).
276 This could simply reflect differences in the patient cohorts or be related to the pathogenesis
277 of the disease, for example due to differences in metabolism or variable death rates across
278 subpopulations of neurons in AD. Studies with larger cohorts will be required to validate and
279 explain this observation.

280

281 To extend these analyses to another tissue not amenable to standard sequencing methods, we
282 studied smooth muscle. Visceral smooth muscle cells are believed to divide infrequently in
283 normal conditions³². We used laser microdissection of histological sections of bladder and
284 colon to collect smooth muscle from 10 donors (**Supplementary Table 1, 2, Extended Data**
285 **Figs. 6b and 10a**). As expected for a polyclonal tissue, standard whole-genome sequencing
286 detected few mutations and at low allele frequencies in these samples (**Extended Data Fig.**
287 **10b, c, Methods**). In contrast, NanoSeq revealed that the substitution and indel burdens
288 increase linearly with age, with ~ 20.7 substitutions per year per diploid genome
289 (CI95%:13.7-28.0) and ~ 1.3 indels per year (95%:0.4-2.3) (**Fig. 3j,k, Supplementary Table**
290 **8**). Despite their different anatomical origin, smooth muscle cells from the bladder and colon
291 walls showed relatively similar mutation rates. Overall, the mutational spectrum of smooth
292 muscle shared some similarities with that of granulocytes and neurons (**Figs. 3l-n** and **1c**),
293 with all three signatures (A-C) accumulating linearly with age (**Extended Data Fig. 7f**). The
294 spectra also resemble that of skeletal muscle satellite cells, studied by in vitro expansion¹¹
295 (**Supplementary Note 10**).

296

297 Altogether, granulocytes, smooth muscle and neurons showed more limited variation in
298 mutation rate and spectra across individuals than has been observed in epithelia exposed to
299 exogenous mutagens, such as skin³, colon⁶ (**Fig. 2c**), bronchus³³ or bladder^{8,34}. This suggests

300 that the variation in endogenous mutagenesis across individuals is modest, at least in the
301 cohorts studied here.

302

303 **Discussion**

304

305 Building on duplex sequencing and BotSeqS, we have developed a sequencing protocol with
306 error rates in single DNA molecules under 5 errors per billion sites. This rate enables the
307 study of mutation rates and signatures in any human tissue or cell population.

308

309 Most of our current knowledge of somatic mutagenesis is restricted to mitotically-active
310 cells. We have exploited the ability to sequence any cell type to study the mutational
311 landscape of non-dividing cells in mitotically-active and inactive tissues. A remarkable
312 observation that emerges from these data is that somatic mutation rates vary modestly (~2-3
313 fold) across a diverse range of somatic cell types, largely independently of cell division rates
314 (**Fig. 3o, p, Supplementary Note 9**). Indeed, similar mutation rates are found in non-
315 dividing cortical neurons, in smooth muscle and in blood; or in colonic epithelium, which
316 divides every few days, and in mostly quiescent hepatocytes¹⁰ or urothelial cells (**Fig. 3o, p**).

317

318 DNA replication and cell division have long been assumed to be major sources of somatic
319 mutations, either due to DNA polymerase errors or the fixation of unrepaired damage during
320 replication³⁵. However, the linear accumulation of somatic mutations in post-mitotic neurons,
321 with similar rates and signatures to some mitotically-active tissues, indicates that dominant
322 mutational processes can occur independently of cell division. These mutations may result
323 from the interplay between endogenous DNA damage and repair that cells are engaged in at
324 all times. The similar mutation burden and signatures in granulocytes and in the stem cells
325 responsible for long-term maintenance of blood, despite a different divisional load, could also
326 be consistent with a time-dependent rather than a division-dependent accumulation of
327 somatic mutations during haematopoiesis. Altogether, division-independent mutational
328 processes may play a larger role in adult mutagenesis than it is commonly assumed.

329

330 In addition to enabling studies on somatic mutagenesis in any tissue, the ability to accurately
331 detect mutations in single molecules of DNA has wider applications. NanoSeq could be used
332 for mutagenesis screens and in vitro studies, exposing cell cultures or experimental models to
333 different mutagens and quantifying mutagenesis across the genome and over time, without
334 the need of single-cell bottlenecks^{36,37}. Sonication followed by exonuclease digestion opens
335 the door to targeted applications, to study the landscape of driver or pathogenic mutations in
336 polyclonal samples, across tissues and conditions. Being insensitive to clonality, NanoSeq
337 can also be used to efficiently and accurately quantify somatic mutation rates and signatures
338 in non-invasive tissue samples, enabling studies of somatic mutagenesis in large-scale
339 cohorts, across genetic backgrounds, exposures and risk factors, in health and disease.

340

341 **References**

342

343

- 344 1 Kennedy, S. R., Loeb, L. A. & Herr, A. J. Somatic mutations in aging, cancer and
345 neurodegeneration. *Mech Ageing Dev* **133**, 118-126, doi:10.1016/j.mad.2011.10.009
346 (2012).
- 347 2 Vogelstein, B. *et al.* Cancer genome landscapes. *Science* **339**, 1546-1558,
348 doi:10.1126/science.1235122 (2013).

349 3 Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive selection
350 of somatic mutations in normal human skin. *Science* **348**, 880-886,
351 doi:10.1126/science.aaa6806 (2015).

352 4 Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age.
353 *Science* **362**, 911-917, doi:10.1126/science.aau3879 (2018).

354 5 Yizhak, K. *et al.* RNA sequence analysis reveals macroscopic somatic clonal
355 expansion across normal tissues. *Science* **364**, doi:10.1126/science.aaw0726 (2019).

356 6 Lee-Six, H. *et al.* The landscape of somatic mutation in normal colorectal epithelial
357 cells. *Nature* **574**, 532-537, doi:10.1038/s41586-019-1672-7 (2019).

358 7 Brunner, S. F. *et al.* Somatic mutations and clonal dynamics in healthy and cirrhotic
359 human liver. *Nature* **574**, 538-542, doi:10.1038/s41586-019-1670-9 (2019).

360 8 Li, R. *et al.* Macroscopic somatic clonal expansion in morphologically normal human
361 urothelium. *Science* **370**, 82-89, doi:10.1126/science.aba7300 (2020).

362 9 Welch, J. S. *et al.* The origin and evolution of mutations in acute myeloid leukemia.
363 *Cell* **150**, 264-278, doi:10.1016/j.cell.2012.06.023 (2012).

364 10 Blokzijl, F. *et al.* Tissue-specific mutation accumulation in human adult stem cells
365 during life. *Nature* **538**, 260-264, doi:10.1038/nature19768 (2016).

366 11 Franco, I. *et al.* Somatic mutagenesis in satellite cells associates with human skeletal
367 muscle aging. *Nat Commun* **9**, 800, doi:10.1038/s41467-018-03244-6 (2018).

368 12 Lodato, M. A. *et al.* Somatic mutation in single human neurons tracks developmental
369 and transcriptional history. *Science* **350**, 94-98, doi:10.1126/science.aab1785 (2015).

370 13 Lodato, M. A. *et al.* Aging and neurodegeneration are associated with increased
371 mutations in single human neurons. *Science* **359**, 555-559,
372 doi:10.1126/science.aao4426 (2018).

373 14 Brazhnik, K. *et al.* Single-cell analysis reveals different age-related somatic mutation
374 profiles between stem and differentiated cells in human liver. *Sci Adv* **6**, eaax2659,
375 doi:10.1126/sciadv.aax2659 (2020).

376 15 Xing, D., Tan, L., Chang, C. H., Li, H. & Xie, X. S. Accurate SNV detection in single
377 cells by transposon-based whole-genome amplification of complementary strands.
378 *Proc Natl Acad Sci U S A* **118**, doi:10.1073/pnas.2013106118 (2021).

379 16 Petljak, M. *et al.* Characterizing Mutational Signatures in Human Cancer Cell Lines
380 Reveals Episodic APOBEC Mutagenesis. *Cell* **176**, 1282-1294.e1220,
381 doi:10.1016/j.cell.2019.02.012 (2019).

382 17 Salk, J. J., Schmitt, M. W. & Loeb, L. A. Enhancing the accuracy of next-generation
383 sequencing for detecting rare and subclonal mutations. *Nat Rev Genet* **19**, 269-285,
384 doi:10.1038/nrg.2017.117 (2018).

385 18 Schmitt, M. W. *et al.* Detection of ultra-rare mutations by next-generation
386 sequencing. *Proc Natl Acad Sci U S A* **109**, 14508-14513,
387 doi:10.1073/pnas.1208715109 (2012).

388 19 Kennedy, S. R. *et al.* Detecting ultralow-frequency mutations by Duplex Sequencing.
389 *Nat Protoc* **9**, 2586-2606, doi:10.1038/nprot.2014.170 (2014).

390 20 Hoang, M. L. *et al.* Genome-wide quantification of rare somatic mutations in normal
391 human tissues using massively parallel sequencing. *Proc Natl Acad Sci U S A* **113**,
392 9846-9851, doi:10.1073/pnas.1607794113 (2016).

393 21 You, X. *et al.* Detection of genome-wide low-frequency mutations with Paired-End
394 and Complementary Consensus Sequencing (PECC-Seq) revealed end-repair-derived
395 artifacts as residual errors. *Arch Toxicol* **94**, 3475-3485, doi:10.1007/s00204-020-
396 02832-0 (2020).

- 397 22 Costello, M. *et al.* Discovery and characterization of artifactual mutations in deep
398 coverage targeted capture sequencing data due to oxidative DNA damage during
399 sample preparation. *Nucleic Acids Res* **41**, e67, doi:10.1093/nar/gks1443 (2013).
- 400 23 Kong, A. *et al.* Rate of de novo mutations and the importance of father's age to
401 disease risk. *Nature* **488**, 471-475, doi:10.1038/nature11396 (2012).
- 402 24 Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat Genet*
403 **48**, 126-133, doi:10.1038/ng.3469 (2016).
- 404 25 Wyles, S. P., Brandt, E. B. & Nelson, T. J. Stem cells: the pursuit of genomic
405 stability. *Int J Mol Sci* **15**, 20948-20967, doi:10.3390/ijms151120948 (2014).
- 406 26 Lee-Six, H. *et al.* Population dynamics of normal human blood inferred from somatic
407 mutations. *Nature* **561**, 473-478, doi:10.1038/s41586-018-0497-0 (2018).
- 408 27 Nicholson, A. M. *et al.* Fixation and Spread of Somatic Mutations in Adult Human
409 Colonic Epithelium. *Cell Stem Cell* **22**, 909-918.e908,
410 doi:10.1016/j.stem.2018.04.020 (2018).
- 411 28 Pleguezuelos-Manzano, C. *et al.* Mutational signature in colorectal cancer caused by
412 genotoxic pks(+) *E. coli*. *Nature* **580**, 269-273, doi:10.1038/s41586-020-2080-8
413 (2020).
- 414 29 Poduri, A., Evrony, G. D., Cai, X. & Walsh, C. A. Somatic mutation, genomic
415 variation, and neurological disease. *Science* **341**, 1237758,
416 doi:10.1126/science.1237758 (2013).
- 417 30 Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer.
418 *Nature* **578**, 94-101, doi:10.1038/s41586-020-1943-3 (2020).
- 419 31 Rheinbay, E. *et al.* Analyses of non-coding somatic drivers in 2,658 cancer whole
420 genomes. *Nature* **578**, 102-111, doi:10.1038/s41586-020-1965-x (2020).
- 421 32 Gabella, G. Cells of visceral smooth muscles. *J Smooth Muscle Res* **48**, 65-95,
422 doi:10.1540/jsmr.48.65 (2012).
- 423 33 Yoshida, K. *et al.* Tobacco smoking and somatic mutations in human bronchial
424 epithelium. *Nature* **578**, 266-272, doi:10.1038/s41586-020-1961-1 (2020).
- 425 34 Lawson, A. R. J. *et al.* Extensive heterogeneity in somatic mutation and selection in
426 the human bladder. *Science* **370**, 75-82, doi:10.1126/science.aba8347 (2020).
- 427 35 Gao, Z., Wyman, M. J., Sella, G. & Przeworski, M. Interpreting the Dependence of
428 Mutation Rates on Age and Time. *PLoS Biol* **14**, e1002355,
429 doi:10.1371/journal.pbio.1002355 (2016).
- 430 36 Kucab, J. E. *et al.* A Compendium of Mutational Signatures of Environmental Agents.
431 *Cell* **177**, 821-836.e816, doi:10.1016/j.cell.2019.03.001 (2019).
- 432 37 Matsumura, S. *et al.* Genome-wide somatic mutation analysis via Hawk-Seq™
433 reveals mutation profiles associated with chemical mutagens. *Arch Toxicol* **93**, 2689-
434 2701, doi:10.1007/s00204-019-02541-3 (2019).

435

436 **Figure legends**

437

438 **Figure 1 | BotSeqS and NanoSeq sequencing protocols.** **a**, Duplex sequencing protocol. **b**,
439 BotSeqS mutation burden estimates in cord blood granulocytes compared to 100 single-cell
440 derived blood colonies from two donors. **c**, BotSeqS and NanoSeq substitution profiles for
441 cord blood granulocytes, and cosine similarities (**Methods**) with the cord blood colonies
442 profile. **d**, Substitution imbalances are present in standard BotSeqS protocols but absent from
443 NanoSeq (**Extended Data Figs. 1 and 2** show further details for a library of granulocytes
444 from a 59-year-old donor). Imbalances were tested with a binomial test (p of 0.5) and p -

445 values were corrected with Benjamini and Hochberg's FDR method. **e**, Standard BotSeqS
446 (top) and NanoSeq protocols (bottom) for library preparation. **f**, **g**, NanoSeq mutation burden
447 estimates for cord blood granulocytes (S1/PD48442, n=6 libraries; S2/PD47269, n=1) and
448 sperm from a 21-year-old donor (n=7) compared to blood colonies and sperm estimates,
449 respectively; **b**, **f**, **g**, Bars show point estimates and their 95% Poisson confidence intervals. **b**,
450 **f**, Box plot shows the interquartile range, median, 95% confidence interval for the median,
451 and outliers as grey dots. **f**, The mean and its 95% confidence interval are shown in red.

452

453

454 **Figure 2 | Mutation in stem and differentiated cells.** **a**, Schematic representation of the
455 hematopoietic lineage showing which cell types and donors were analysed. **b**, Substitutions
456 per cell for donors of different ages, comparing granulocytes and single-cell derived blood
457 colonies. NanoSeq estimates for granulocytes (red dots) obtained for one library per donor
458 except for donors of ages 54 (n=2), 63 (n=2), and 59 (n=5). Standard sequencing estimates
459 are shown as box plots and based on 10 colonies per donor, except for the 59-year-old
460 (n=110) and cord blood (n=100). **c**, Granulocytes and blood colonies substitution profiles and
461 their cosine similarity (**Methods**) for the 59-year-old donor. **d**, Burden estimates in colonic
462 crypts from three donors, comparing standard methods (box plots) and NanoSeq (red dots;
463 n=3, 2 and 2 libraries per donor); **e**, Accumulation of substitutions throughout life in colonic
464 crypts from five donors, excluding substitutions attributed to the episodic colibactin
465 signature. **b**, **d**, Dots and lines show point estimates and their corresponding 95% Poisson
466 confidence intervals, respectively. **b**, **d**, Box plots show the interquartile range, median, 95%
467 confidence interval for the median, with outliers as grey dots. **b**, **e**, Linear mixed regression
468 models for granulocytes (red dashed line), blood colonies (dark cyan), and colonic crypts
469 (black), with 95% confidence intervals calculated through parametric bootstrapping
470 (**Methods**). Regression intercepts and slopes are provided in **Supplementary Table 8**.

471

472 **Figure 3 | Mutation landscape in neurons and smooth muscle.** **a**, **b**, Accumulation of
473 substitutions and indels in neurons throughout life, including healthy (n=8) and Alzheimer's
474 disease (n=9) donors. **c**, **d**, Substitution and indel spectra in neurons; a description of each
475 type of indel can be found in **Extended Data Fig. 5d**. **e**, Signature decomposition. **f**, **g**, Indel
476 and substitution rates in genes in the whole cohort by level of expression. **h**, Substitution
477 rates in transcribed and quiescent/heterochromatin DNA across different cell types (spectra in
478 **Extended Data Fig. 8a**). **i**, Contribution of signatures A, B and C in neurons. **j**, **k**,
479 Substitutions and indels per cell in smooth muscle from 10 donors spanning different ages
480 (n=2 libraries for donors aged 54 and 68). **l**, **m**, Substitution and indel spectra in smooth
481 muscle. **n**, Exposure to signatures A, B and C in smooth muscle samples. **o**, **p**, Substitution
482 and indel accumulation per year across different cell types; 95% confidence intervals
483 estimated through simple (neurons and urothelium) or mixed effect (rest) linear regression
484 with intercept=0; vertical lines show regression 95% confidence intervals. **a**, **b**, **f**, **g**, **h**, **j**, **k**,
485 vertical lines show Poisson 95% confidence intervals. **a**, **b**, Linear regression model as
486 dashed lines, showing 95% confidence interval as grey areas. **j**, **k**, Linear mixed effect
487 regressions as dashed lines, showing 95% confidence intervals obtained through parametric
488 bootstrapping (**Methods**) as grey areas. Regression results with free or zero intercept are
489 provided in **Supplementary Table 8**.

490

491

492 **Methods**

493

494 **Sample collection and ethics**

495

496 All samples were collected with informed consent from all human research participants or
497 their families. The haematological samples in the study were obtained from the Cambridge
498 Blood and Stem Cell Biobank, the Cambridge Biorepository for Translational Medicine, and
499 the Cambridge Bioresource (REC references: 07-MRE05-44, 18/EE/0199, 15/EE/0152 -
500 NRES Committee East of England - Cambridge South). Sperm samples were collected under
501 REC ethics approval EC04/015, London - Westminster REC; 16/NE/003, NRES Committee
502 North East-Newcastle and North Tyneside 1. Colon and bladder tissue were collected by the
503 Cambridge Biorepository for Translational Medicine (REC reference: 15/EE/0152 NRES
504 Committee East of England – Cambridge South). Frozen biopsies of frontal cortex from
505 healthy and Alzheimer’s disease donors were collected by the Cambridge Brain Bank
506 (**Supplementary Tables 1, 2**; REC ethics approval: 10/H0308/56, East of England,
507 Nottingham).

508

509 **Granulocytes and HSC/MPP colonies: sorting, colony growth and mutation calling**

510

511 We use two different terms to refer to colonies derived from haematopoietic stem cells (HSC)
512 or progenitor cells, depending on the membrane markers used for cell sorting: HSPCs, which
513 refer to CD34+ pools, and HSC/MPPs, which refer to CD34+ CD38- CD45RA- cells.

514

515 A sample of granulocytes from a 59-year-old male donor (PD43976_59yo) from whom 110
516 HSPC colonies were available²⁶ was used for initial validation of the BotSeqS and NanoSeq
517 protocols (**Supplementary Tables 1, 2**). To estimate the NanoSeq error rate, cord blood
518 granulocytes from two neonatal donors were sequenced by NanoSeq and the mutation
519 burdens and spectra compared to those from 50 HSC/MPP colonies per donor. For the
520 comparison of differentiated and stem cells, NanoSeq data from granulocytes from 9 donors
521 of different ages were compared to standard sequencing of single-cell derived HSC/MPP
522 colonies from 6 donors (10 HSC/MPP colonies per donor) and 110 HSPC colonies already
523 available from a 59-year-old donor²⁶. These 110 HSPC included 67 HSC/MPPs, 32
524 megakaryocyte–erythrocyte progenitors (MEP), 7 granulocyte–macrophage progenitors
525 (GMP) and 4 common myeloid progenitors (CMP).

526

527 For PD43976_59yo, HSPC colonies were grown and mutations called as described in Lee Six
528 *et al.*²⁶. For the remaining donors, whole blood was diluted with PBS and mononuclear cells
529 (MNC) were isolated using lymphoprepTM (STEMCELL Technologies) density gradient
530 centrifugation. The MNC fraction was then removed to a fresh tube, leaving behind the red
531 cell pellet, which also contained the granulocyte fraction. The MNC fraction was depleted of
532 red blood cells by a single 15 min incubation with RBC lysis buffer (BioLegend) at 4°C.
533 Granulocytes were purified from the red cell pellets using 3 incubations (for 20 mins/10
534 mins/10 mins respectively) with RBC lysis buffer (BioLegend) at room temperature. CD34+
535 selection of peripheral blood and cord blood samples was undertaken using the EasySep
536 human whole blood CD34 positive selection kit (STEMCELL Technologies) as per the
537 manufacturer’s instructions. Bone marrow samples did not undergo CD34+ selection prior to
538 sorting.

539

540 MNC or CD34 enriched samples were centrifuged and resuspended in PBS/3%FBS
541 containing an antibody panel consisting of (antibody/fluorochrome): CD3/FITC (1:500),
542 CD90/PE (1:50), CD49f/PECy5 (1:100), CD38/PECy7 (1:100), CD19/A700 (1:300),
543 CD34/APC Cy7 (1:100), CD45RA/BV421 (1:100), and Zombie/Aqua (1:2000).

544

545 Cells were stained (30 minutes at 4°C) in the dark before washing, centrifugation (500 x g at
546 room temperature) and resuspension in PBS/3%FBS for cell sorting. Index sorting of
547 ‘HSC/MPP pool’ cells was performed on a BD AriaIII Cell Sorter (BD Biosciences) at the
548 NIHR Cambridge BRC Cell Phenotyping Hub, as per the gating structure in **Extended Data**
549 **Fig. 6a** (CD34+, CD38- and CD45RA-).

550

551 ‘HSC/MPP pool’ cells were single-cell sorted into Nunc 96 well flat-bottomed TC plates
552 (ThermoFisher) containing 100 µl supplemented StemPro media (Stem Cell Technologies).
553 MEM media contained StemPro Nutrients (0.035%, Stem Cell Technologies), L-Glutamine
554 (1%, ThermoFisher), Penicillin-Streptomycin (1%, ThermoFisher) and cytokines (SCF, 100
555 ng/ml; FLT3, 20 ng/ml; TPO, 100 ng/ml; EPO 3 ng/ml; IL-6, 50 ng/ml; IL-3, 10 ng/ml; IL-
556 11, 50 ng/ml; GM-CSF, 20 ng/ml; IL-2 10 ng/ml; IL-7 20 ng/ml; lipids 50 ng/ml) to promote
557 differentiation towards Myeloid/Erythroid/Megakaryocyte (MEM) and NK lineages. Manual
558 assessment of colony growth was made at 14 days. Colonies were topped up with an
559 additional 50 µL MEM media on day 15 if the colony was $\geq 1/4$ size of well. Following 21 \pm
560 2 days in culture, colonies were selected by size criteria. Colonies ≥ 3000 cells in size were
561 harvested into a U bottomed 96 well plate (ThermoFisher). Plates were then centrifuged (500
562 x g for 5 minutes), media was discarded, and the cells were resuspended in 50 µl PBS prior to
563 freezing at -80°C. Colonies < 3000 cells but > 200 cells in size were harvested into 96 well
564 skirted LoBind plates (Eppendorf) and centrifuged (800 x g for 5 min). Supernatant was
565 removed to 5-10 µL using an aspirator prior to DNA extraction on the fresh cell pellet.

566

567 DNA extraction was performed using the DNeasy 96 blood and tissue plate kit (Qiagen) for
568 larger HSC colonies, or the Arcturus Picopure DNA Extraction kit (ThermoFisher) for
569 smaller HSC colonies. Both kits were used as per the manufacturer’s instructions. Extracted
570 DNA (1-5ng) from each colony was processed using a recently developed low-input
571 enzymatic fragmentation-based library preparation method³⁸. All samples were subjected to
572 whole genome sequencing at 8-35X coverage on either the HiSeq X or the NovaSeq
573 platforms (Illumina) to generate 150 bp paired-end reads. BWA *mem* was used to align
574 sequences to the human reference genome (NCBI build37).

575

576 **Sperm samples**

577

578 DNA was extracted from sperm samples from two donors, aged 21 and 73 years, and
579 sequenced using the NanoSeq protocol. Because of the low mutation burden of the germline,
580 we sequenced 7 separate aliquots of sperm DNA from the 21-year-old donor to estimate the
581 error rate of the NanoSeq protocol (**Supplementary Tables 1, 2**).

582

583 **Laser microdissection of colonic crypts and bladder/colon smooth muscle**

584

585 Colon and bladder biopsies were obtained from deceased organ donors (ranging in age from
586 25 to 78; **Supplementary Table 1**) at the time of organ donation. Different microbiopsies
587 from these specimens have been used in previously published studies^{6,34,39}.

588

589 Colon biopsies were fresh frozen at the time of collection and stored at -80 °C. The colon
590 biopsies subsequently underwent formalin-free fixation for 24 hours in PAXgene Tissue Fix
591 containers (PreAnalytiX, Hombrechtikon, Switzerland) before being transferred to PAXgene
592 STABILIZER solution (PreAnalytiX). Bladder biopsies underwent formalin-free fixation at
593 the time of collection and were stored at -20 °C³⁸.

594

595 Prior to laser-capture microdissection, samples were processed, embedded in paraffin and
596 sectioned as described previously³⁴. Microbiopsies were dissected using an LMD7
597 microscope (Leica Microsystems). Examples of microdissected regions for both specimen
598 types can be found in **Extended Data Figs. 6 and 10**. Proteolysis of isolated regions was
599 performed using an Arcturus PicoPure DNA Extraction Kit (Thermo Fisher Scientific,
600 Waltham, MA, USA). Cell lysate was stored at -20 °C prior to library preparation.

601

602 **Neuron nuclei sorting from frontal cortex samples**

603

604 Neuronal nuclei were isolated, stained and extracted from the frontal cortex samples as per
605 Krishnaswami et al.⁴⁰ using frozen biopsies of frontal cortex from eight healthy and nine
606 Alzheimer's disease donors. Briefly, small cuts of 1-2 mm were taken from fresh frozen
607 samples. Dounce homogenisation was then used to free nuclei before filtration, density
608 centrifugation and immunostaining. Samples were stained using DAPI (Thermo Fisher,
609 D1306) and Milli-Mark™ Anti-NeuN-PE Antibody (1:500; MilliPore, FCMAB317PE). The
610 immunostained samples were then sorted using FACS as per the gating strategy in **Extended**
611 **data Fig. 7a**. 15,000 nuclei were collected into 20 µl Arcturus PicoPure DNA Extraction Kit
612 (Thermo Fisher Scientific) before undergoing digestion. Nuclear lysate was then stored at -
613 20°C prior to library preparation.

614

615 The distributions of NeuN-PE intensities in most samples revealed a bimodal distribution. As
616 a quality control, we fitted a mixture of two Gamma distributions to the NeuN-PE intensities
617 for every samples. Only samples with 10-fold (1 log₁₀ unit) separation between the mean of
618 both peaks were considered for analysis, which led to the exclusion of an outlier sample.

619

620 **BotSeqS and NanoSeq library preparation protocols**

621

622 BotSeqS libraries shown in Fig. 1 and Extended Data Fig. 3 were prepared as follows: DNA
623 was sheared to 450 bp using a Covaris. DNA was cleaned up using a 2.5X Ampure XP
624 (Beckman Coulter) bead ratio. DNA was eluted in 12 µL NFW. 10 µL of the elution product
625 were taken into the ligation reaction consisting in addition of 3.74 µL NEBuffer 4, 3.74 µL
626 10 mM ATP, 0.33 µL xGen Duplex Seq Adapters (IDT 1080799), 0.56 µL T4 DNA ligase
627 (NEB M0202L) and 19.03 µL NFW. The reaction was incubated at 20 °C for 20 min. The
628 DNA was cleaned-up using 37.4 µL Ampure XP beads and DNA was eluted in 50 µL NFW.
629 Libraries were quantified (qPCR) and amplified following the NanoSeq protocol. For the
630 BotSeqS data on granulocytes from a 59-year-old donor, shown in Extended Data Figs. 1 and
631 2, we used an earlier implementation of the protocol. 10 ng of sonicated DNA was end-
632 repaired and ligated using the NEBNext Ultra II kit (New England Biolabs) including 0.66 µl
633 1.5 µM xGen Duplex Seq Adapters - Tech Access (Integrated DNA Technologies, IDT:
634 1080799).

635

636 NanoSeq libraries were prepared as follows: 10 ng of genomic DNA or LCM cut sections in
637 20 µl buffer were purified using 100 µl of a 50:50 water and AMPure XP bead mixture and
638 eluted in 20 µl nuclease free water. 20 µl of the bead suspension was taken forward into an
639 on-bead fragmentation reaction. Fragmentation occurred in a final volume of 25 µl including
640 2.5 µl 10x CutSmart buffer (500 mM Potassium Acetate, 200 mM Tris-acetate, 100 mM
641 Magnesium Acetate, 1 mg/ml BSA, pH 7.9 at 25°C), 0.5 µl 5 U/µl HpyCH4V
642 (**Supplementary Note 2**), and 2 µl NFW. Fragmentation reactions were incubated at 37 °C
643 for 15 min, purified with 2.5x AMPure XP beads and resuspended in 15 µl nuclease-free

644 water. Fragmented DNA was A-tailed in 15 μ l reactions including 10 μ l fragmentation
645 product, 1.5 μ l 10x NEBuffer 4 (500 mM Potassium Acetate, 200 mM Tris-acetate, 100 mM
646 Magnesium Acetate, 10 mM DTT, pH 7.9 at 25°C), 0.15 μ l 5 U/ μ l Klenow fragment (3'→5'
647 exo-, New England Biolabs), either 1.5 μ l 1 mM dATP or 1.5 μ l 1 mM equimolar
648 dATP/ddBTPs (**Supplementary Note 3**), and 1.85 μ l NFW. Reactions were incubated at 37
649 °C for 30 mins. The 15 μ l A-tailing reaction product was added to 22.4 μ l ligation mix, which
650 consisted of 2.24 μ l 10x NEBuffer 4, 3.74 μ l 10 mM ATP, 0.33 μ l 15 μ M xGen Duplex Seq
651 Adapters (IDT: 1080799), 0.56 μ l 400 U/ μ l T4 DNA ligase (New England Biolabs), and
652 15.53 μ l NFW. Reactions were incubated at 20 °C for 20 min and subsequently purified with
653 1x AMPure XP beads and resuspended in 50 μ l of nuclease free water.

654
655 Mung Bean NanoSeq libraries were prepared as follows: DNA was sheared to an average
656 size of 450 bp using focused ultrasonication (Covaris 644 LE220). Sheared DNA was
657 quantified and 50 ng were used as input per reaction. Mung Bean nuclease (NEB: M0250S)
658 was diluted to 1U, 0.5U or 0.25U/ μ l in 1X Mung Bean nuclease buffer. The Mung Bean
659 reaction was carried out in a final volume of 30 μ l including 2.9 μ l 10X Mung Bean
660 nuclease buffer, 1 μ l diluted Mung Bean nuclease, 10 μ l DNA and 16.1 μ l NFW. The
661 reaction was incubated at 30 °C for 30 min. Then, 1 μ l 0.3% SDS was added and the reaction
662 was cleaned up using 77.5 μ l Ampure XP beads. Samples were eluted in 12 μ l NFW. 10 μ l
663 was used as input into a phosphorylation reaction by adding 1.5 μ l NEBuffer 4 (NEB
664 B7004S), 1.5 μ l 10 mM ATP (Fisher Scientific 10304340), 0.6 μ l T4 Polynucleotide
665 Kinase (NEB M0201S) and 1.4 μ l NFW. The reaction was incubated at 37 °C for 30 min. 13
666 μ l were taken forward into an A-tailing reaction, adding 0.2 μ l NEBuffer 4, 1.5 μ l 1 mM
667 dATP/ddBTP (NEB N0440S/GE Healthcare 27204501), 0.15 μ l Klenow fragment (3'→5'
668 exo-, NEB M0212L) and 0.15 μ l NFW. The reaction was incubated at 37 °C for 30 min. The
669 whole 15 μ l were taken into the ligation reaction mix, which consisted of 2.24 μ l NEBuffer
670 4, 3.74 μ l 10 mM ATP, 0.33 μ l xGen Duplex Seq Adapters (IDT 1080799), 0.56 μ l T4
671 DNA ligase (NEB M0202L) and 15.53 μ l NFW. The reaction was incubated at 20 °C for 20
672 min. The DNA was cleaned-up using 37.4 μ l Ampure XP beads and DNA was eluted in 50
673 μ l NFW. Libraries were quantified and amplified following the NanoSeq protocol.

674 675 **DNA quantification, dilution and PCR amplification**

676
677 DNA was quantified by qPCR using a KAPA library quantification kit (KK4835). The
678 supplied primer premix was first added to the supplied KAPA SYBR FAST master mix. In
679 addition, 20 μ l of 100 μ M NanoqPCR1 primer (HPLC: 5'-ACACTCTTTCCTACACGAC-
680 3') and 20 μ l of 100 μ M NanoqPCR2 primer (HPLC: 5'-GTGACTGGAGTTCAGACGTG-
681 3') were added to the KAPA SYBR FAST master mix. Samples were diluted 1 in 500 using
682 nuclease-free water and reactions were set up in a 10 μ l reaction volume (6 μ l master mix, 2
683 μ l sample/standard, 2 μ l water) in a 384 well plate. Samples were run on the Roche 480
684 Lightcycler and analysed using absolute quantification (2nd Derivative Maximum Method)
685 with the high sensitivity algorithm. nM (fmol/ μ l) was determined as follows: mean of sample
686 concentration x dilution factor (500) x 452/573/1000 (where 452 is the size of the standard in
687 bp and 573 is the proxy for the average fragment length of the library in bp), and multiplied
688 by an adjustment factor of 1.5. Samples were diluted to the desired fmol amount (typically
689 0.3 fmol for a 15x run) in 25 μ l using nuclease free water.

690
691 Libraries were subsequently PCR amplified in a 50 μ l reaction volume comprising of 25 μ l
692 sample, 25 μ l NEBNext Ultra II Q5 Master Mix and UDI containing PCR primers (dried).
693 The reaction was cycled as follows: step1: 98 °C 30 seconds, step2: 98 °C 10 seconds, step3:

694 65 °C 75 seconds, step4: return to step2 13 times, step5: 65 °C for 5 min, step6: hold at 4 °C.
695 The number of PCR cycles is dependent upon the input: 0.1 fmol = 16 cycles, 0.3 fmol = 14
696 cycles, 0.6 fmol = 13 cycles, 5 fmol = 10 cycles.

697

698 The PCR product was subsequently cleaned up using two consecutive 0.7x AMPure XP
699 clean-ups. Each sample was quantified using the AccuClear Ultra High Sensitivity dsDNA
700 Quantification kit (Biotium) and pooled. Libraries were sequenced on Illumina sequencing
701 platforms e.g. NovaSeq using 150 paired-end reads.

702

703 **Library dilution and sequencing efficiency**

704

705 The efficiency and cost-effectiveness of duplex sequencing depends on optimising the
706 duplicate rate to maximise the number of read bundles (defined as a family of PCR
707 duplicates) with at least 2 duplicate reads from each original strand. Too high duplicate rates
708 result in few read bundles of unnecessarily large sizes, whereas too low duplicate rates result
709 in many read bundles with few having two or more read pairs from each strand.

710

711 To maximise the efficiency of the protocol, we studied analytically and empirically the
712 relationship between the number of DNA molecules in the library (library complexity) and
713 the resulting duplicate rate as a function of the number of read pairs sequenced. We found
714 that optimal duplicate rates and optimal efficiency can be ensured across a wide range of
715 samples. If we assume negligible PCR biases, with copies from all original ligated DNA
716 fragments represented in equimolar amounts in the amplified library, the bundle size
717 distribution of observed reads can be modelled as a zero-truncated Poisson distribution. Let r
718 (sequence ratio) be the ratio between the number of sequenced reads and the number of
719 amplifiable DNA fragments in the original library. The mean read bundle size (m) can then
720 be estimated as the mean of the zero-truncated Poisson distribution: $m = \frac{r}{1-e^{-r}}$. This
721 parameter then enables a simple estimation of the duplicate rate of a library (d , defined as the
722 fraction of reads that are duplicate copies, and identified as reads having the same barcodes
723 and the same 5' and 3' coordinates): $d = \frac{m-1}{m} = 1 - \frac{1}{m} = 1 - \frac{1-e^{-r}}{r}$.

724

725 We can define the efficiency of a duplex sequencing library (E) as the ratio between the
726 number of base pairs with duplex coverage (bundles with ≥ 2 reads from both strands) and

727 the number of base pairs sequenced. This can be modelled as: $E = \frac{P(x \geq 2; \frac{r}{2})^2}{m}$, where the

728 numerator is the probability of a read bundle having at least two reads from both strands (i.e.
729 usable bundles), based on the zero-truncated Poisson distribution (denoted as P), and the

730 denominator is the sequence investment in each read bundle (i.e. the average read bundle
731 size). Based on this equation, we can estimate numerically that the optimal duplicate rate is

732 ~81% (**Extended Data Fig. 4a, Supplementary Code**) and that duplicate rates between 65-
733 90% would yield $\geq 80\%$ of the maximum attainable efficiency. In terms of r , the optimum r is

734 5.1 read pairs sequenced per original DNA fragment (r_{opt}), with values within 2.7-9.6
735 yielding $\geq 80\%$ of the maximum efficiency. Knowing the concentration of a NanoSeq (or

736 BotSeqS) library in fmol/ μ l (estimated using a qPCR reaction on an aliquot of the library),
737 we can use r_{opt} to calculate the volume of library that needs to be amplified to yield optimal

738 duplicate rates (i.e. maximum duplex efficiency), as a function of the desired amount of raw
739 sequencing: $fmol_{opt} = \frac{N}{f r_{opt}}$. Here, N is the number of paired-end reads that will be
740 sequenced and f is the number of DNA fragments per fmol of library (referring specifically to

741 ligated and amplifiable fragments within the size selection range). Using an initial set of
742 libraries, we compared a range of library inputs (fmol) to the estimated number of unique
743 molecules in the library inferred from the sequencing data (using Piccard's software). This
744 analysis revealed that, for our choice of restriction enzyme and size selection conditions, f
745 approximately equated to 10^8 fragments/fmol (**Supplementary Code**).

746
747 Using the above equation, we can optimise the efficiency of NanoSeq independently of the
748 input amount of DNA in a given sample. For example, ~ 0.3 fmols of library yield optimal
749 duplicate rates when using 150 million 150 bp paired-end reads, which are the equivalent of
750 $\sim 15x$ coverage in standard human whole-genome sequencing. ~ 0.6 fmol yield optimal
751 efficiency when using 300 million reads (30x whole-genome equivalent). Note that, as
752 predicted by the equations above, deviations ~ 2 -fold from r_{opt} still yield high efficiency.
753 Using these equations, we reliably obtained near-optimal duplicate rates from a wide
754 diversity of samples (**Extended Data Fig. 4, Supplementary Table 2**). Overall, we found
755 that $\sim 30x$ of standard sequencing output ($\sim 300 \times 10^6$ 150bp PE reads) yielded approximately
756 3 Gb of high-accuracy duplex coverage (a haploid genome equivalent) after application of all
757 computational filters.

758
759 Our choices of restriction enzyme and size selection restrict the coverage to $\sim 30\%$ of the
760 human genome. Although the covered regions are sufficiently diverse to enable unbiased
761 estimates of burden and signatures (**Methods**), applications that require full genome
762 coverage, such as targeted sequencing, would require alternative fragmentation strategies.
763 One option may be exonuclease blunting after sonication, instead of end repair. Nevertheless,
764 for the study of burden and signatures, the use of restriction enzymes has two interesting
765 advantages. First, this protocol is able to work with very low inputs of DNA. We estimated
766 library yields for a range of input DNA amounts (**Extended Data Fig. 4b**) and found that the
767 minimum DNA input required to obtain 0.3 fmol for a 15x run (corresponding to about 1.5-3
768 Gb of effective duplex coverage) was ~ 1 ng of input DNA. This low-input requirement
769 enables the application of NanoSeq to microscopic areas of tissue (as shown for colonic
770 crypts and smooth muscle) and to rare cell populations using flow sorting. A second
771 advantage is that, since coverage is concentrated in $\sim 30\%$ of the human genome, matched
772 normal samples can be sequenced at lower cost by using undiluted NanoSeq libraries (≥ 3
773 fmol of library sequenced at 8x genome equivalent is enough to provide high matched normal
774 coverage in the 30% of informative genome).

775 776 **Sequencing, pre-processing and filtering of BotSeqS and NanoSeq libraries**

777
778 Standard sequencing matched-normal libraries were aligned to the human reference genome
779 (GRCh37, hs37d5 build) using BWA-MEM v0.7.5a-r405⁴¹ with default parameters.
780 Alignments were sorted by coordinate and read duplicates were marked using biobambam2⁴²
781 v2.076 bamsormadup. Matched-normal reads were filtered if marked as duplicate,
782 supplementary, QC fail, unmapped or secondary alignments. For some samples, as described
783 above, instead of standard whole-genome sequencing, we used undiluted NanoSeq libraries
784 (typically ~ 5 fmol) as matched normals, reducing the costs of sequencing matched normal
785 samples.

786
787 NanoSeq and BotSeqS libraries were sequenced using 150 bp paired-end reads, on
788 HiSeq2500, HiSeqX and NovaSeq platforms.

789

790 NanoSeq sequencing reads begin with adapter sequences: NNNT or NNNXT for BotSeqS
 791 libraries and NNNTCA or NNNXTCA for HpyCH4V libraries (HpyCH4V cuts at TGCA
 792 motifs). NNN is a random three nucleotide barcode, T is the adapter overhang and X is a
 793 ‘spacer’ nucleotide designed to increase nucleotide diversity in the sequencing run. We used
 794 a custom Python script to process demultiplexed fastq files by extracting the three-nucleotide
 795 barcode, clipping remaining adapter bases (2 bases for BotSeqS and 4 bases for NanoSeq
 796 libraries) and appending barcode sequences to the fastq header. Barcodes with non-canonical
 797 bases (not A, C, G or T) were filtered out. Reads were aligned to hs37d5 using *bwa mem*
 798 (v0.7.5a-r405), using the -C option to append barcode sequences to alignments. Alignments
 799 were sorted by coordinate, duplicates were marked, and reads were annotated with read
 800 coordinate, mate coordinate and optical duplicate auxiliary tags using *biobambam2* v2.076
 801 *bamsormadup* and *bammarkduplicatesopt* (optminpixeldif=2500). Reads were filtered when
 802 they were not marked as proper-pairs or were marked as optical duplicate, supplementary,
 803 QC fail, unmapped or secondary alignments. Each read was marked with an auxiliary tag
 804 comprised of reference name, sorted read and mate fragmentation breakpoints, forward and
 805 reverse read barcodes, and read strand.

806

807 **Consensus base quality scores**

808

809 Bayes’ theorem was used to compute the posterior probability of each base call B given the
 810 pileup of reads D from one strand of a template molecule at one genomic position. There are
 811 four possible genotypes $i \in (A, C, G, T)$. The posterior probability is calculated using:

812

$$P(B|D) = \frac{P(B)P(D|B)}{\sum_i P(B_i)P(D|B_i)}$$

813

814 Under a uniform prior, where any of the four possible genotypes are equally likely, the
 815 equation can be simplified to:

816

$$P(B|D) = \frac{P(D|B)}{\sum_i P(D|B_i)}$$

817

818 To calculate $P(D|B)$, information is integrated from reads in D , where $b_j \in (A, C, G, T)$ is the
 819 base of read $j = 1 \dots d$:

820

$$P(D|B_i) = \prod_{j=1}^{j=d} P(b_j|B_i)$$

821

822 To calculate $P(b_j|B_i)$ we use the probability that base b_j is an error, calculated from its Phred
 823 quality score q_j :

824

$$P(b_j|G_i) = 1 - e_j \text{ if } b_j = B_i, \text{ otherwise } e_j/3$$

825

826 where

827

$$e_j = 10^{-\frac{q_j}{10}}$$

828

829 We note that the final probability $P(D|B)$ is the probability that the base call is correct after
830 sequencing and not the probability that the base represents the correct genotype of the
831 original template strand, where independence between observations cannot be assumed.
832 $P(B|D)$ is rescaled into a Phred quality score Q using:

833

$$Q = -10 \log_{10} P(B|D)$$

834

835 In cases where the two read mates overlap, the consensus base quality is calculated using
836 both forward and reverse reads.

837

838 **Base calling and filtering**

839

840 We developed a set of filters that successfully reduced false positive calls. An important
841 feature of the bioinformatic pipeline is that we apply the same filters to call reference and
842 mutated bases, which allows direct calculation of mutation rates.

843

844 The calling method requires a matched normal to filter out germline SNPs. An additional
845 mask to filter sites that are problematic is also advisable. This matched normal can be
846 obtained by standard protocols or by sequencing undiluted NanoSeq libraries (≥ 3 fmol), as
847 explained above.

848

849 The same filters were applied to NanoSeq and BotSeq data. (a) We require that each read
850 bundle (i.e. group of PCR duplicates) has at least two reads from each of the two original
851 DNA strands. (b) The consensus base quality score should be at least 60 (this guarantees that
852 there is strong support for a given base call from the duplicate reads that form a read bundle).
853 (c) the minimum difference between the primary alignment score (AS) and the secondary
854 alignment score (XS) should be higher than 50, to keep only read pairs with unambiguous
855 mapping (for sites where the two mates overlap the minimum of the average AS-XS for
856 forward and reverse mates is taken). This filter is essential to remove mapping artefacts and a
857 minimum AS-XS of 50 is applied also to the matched normal. (d) The average number of
858 mismatches (NM) in a group of reads (forward or reverse) should not be higher than 2, either
859 in the matched normal or the sample at hand. To avoid a bias in the filtering of mutation and
860 reference calls, where a consensus base call is different from the reference, mismatches from
861 that call are not considered when calculating the number of mismatches in the read. For sites
862 where the two mates overlap, the maximum of the average NM for forward and reverse mates
863 is taken. (e) No 5' clips are allowed. (f) No improper pairs are allowed in the read bundle to
864 avoid unreliable mappings. (g) Base calls in read ends, defined as those within 8 bp from the
865 5' or 3' ends, are discarded because these regions are more likely to be unreliably mapped,
866 especially when there are nearby indels. (h) Reads in the read bundle must contain no indels
867 (except for indel calling). (i) The matched normal must have $\geq 15x$ coverage at a given site to
868 make the risk of undetected heterozygous SNPs negligible. For non-neat matched normals we
869 also require that there are at least five reads aligned to each strand. (k) When a mutation is to
870 be called, we require that the base is not seen with a frequency higher than 0.01 in the
871 matched normal. (l) A site should not overlap the common SNP and noisy sites masks (see
872 **Genome masks**). Base calls failing this requirement are also counted to obtain a qualitative
873 diagnostic of potential contamination of the input DNA with DNA from a different
874 individual.

875

876 **Indel calling**

877

878 To call indels we first identify read bundles with potential indels, defined as those containing
879 sites with at least 90% of forward and reverse reads having an indel. Read bundles with AS-
880 XS ≤ 50 , 5' clipping or with coverage in the matched normal lower than 16 were filtered out.
881 Indels close to read ends (10 bp) were not called. For each of the read bundles potentially
882 containing an indel, the corresponding reads were extracted from the BAM file, removing
883 PCR duplicate flags and creating a mini read bundle BAM. For each of the read bundle
884 BAMs we run samtools mpileup to generate genotype likelihoods in BCF format: samtools
885 mpileup --no-BAQ -d 250 -m 2 -F 0.5 -r \$chr:\$start-\$end --BCF --output-tags
886 DP,DV,DP4,SP -f \$ref_genome -o genotype_likelihoods.bcf read_bundle.bam, where \$chr,
887 \$start and \$end are the mapping coordinates of the read bundle. Next, we call indels and
888 normalise the output using the following three bcftools commands: 1) bcftools index -f
889 genotype_likelihoods.bcf genotype_likelihoods.indexed.bcf; 2) bcftools call --skip-variants snps
890 --multiallelic-caller --variants-only -O v genotype_likelihoods.bcf -o bcftools.tmp.vcf; and 3)
891 bcftools norm -f \$ref_genome bcftools.tmp.vcf > bcftools.tmp2.vcf.

892

893 For each of the sites involved in an indel we check whether it overlaps a site masked by our
894 common SNP and noise masks (see **Genome masks**), in which case the indel is flagged as
895 MASKED and not further analysed.

896

897 The final step involves revisiting the matched normal to inspect if there are indels in a
898 window of ± 5 bp around each candidate indel. For this step we use the bam2R function from
899 R package *deepSNV*⁴³. Reads with mapping quality lower than 10 or with any of the
900 following flags are ignored: "read unmapped", "not primary alignment", "read fails
901 platform/vendor quality checks", "read is PCR or optical duplicate", and "supplementary
902 alignment". If the proportion of indels in the matched normal within the ± 5 bp window
903 around the candidate somatic indel is higher than 1%, the indel is disregarded.

904

905 **Substitution imbalances**

906

907 To detect asymmetries in substitution patterns, variants were assigned to the forward or
908 reverse strand according to their distance from fragmentation breakpoints. Variants closest to
909 the 5' of the forward read were assigned to the forward strand. Variants closest to the 5' of the
910 reverse read were assigned to the reverse strand and reverse complemented. Variants
911 equidistant from both fragmentation breakpoints were not counted.

912

913 **Genome masks**

914

915 We applied two masks to filter duplex sequencing data. The first mask comprised common
916 SNPs and spanned a total of 27,204,965 bp. Autosomal and X-chromosome common SNPs
917 were defined as SNPs with allele frequency (AF) $> 0.1\%$ and a "PASS" flag in gnomAD. Y-
918 chromosome and mitochondrial SNPs were defined as SNPs with AF $> 0.1\%$ from 1000
919 Genomes Project (1KGP) data^{44,45}. This SNP mask is important to reduce the impact of
920 potential inter-individual DNA contamination (**Supplementary Note 6**).

921

922 A second mask was developed to remove unreliable calls or sites prone to alignment
923 artefacts. To build this noise mask we gathered together gnomAD indel calls with AF $> 1\%$
924 and SNP calls with AF $> 0.1\%$ that were not flagged as "PASS". The noise mask also contains
925 sites with elevated error-rates. To generate the mask, mismatch rates were calculated for

926 every genomic position across a panel of 448 in-house standard whole-genome samples. Sites
927 with mismatch rates (coverage-weighted mean VAF) > 0.01 were incorporated into the noise
928 mask. Altogether, the second mask comprised 22,474,160 bp.

929

930 Both masks are available at https://github.com/fa8sanger/NanoSeq_Paper_Code.

931

932 **Detection of human DNA contamination**

933

934 Contamination of duplex sequencing libraries with DNA from other individuals could
935 artificially inflate mutation burden estimates, mainly because germline SNPs in the
936 contaminant DNA may appear as somatic mutations.

937

938 Even a small percentage of contamination can have a large impact on burden estimates. The
939 burden associated to SNPs in the contaminant would be:

940

$$Burden_{SNP} = \frac{N_{SNP} * f_{cont}}{G}$$

941

942 being N_{SNP} the number of SNPs in the contaminant not shared with the sample at hand, f_{cont}
943 the contamination fraction and G the size of the diploid human genome. Accordingly, 1%
944 contamination would result in a $Burden_{SNP}$ of $\sim 5 \times 10^{-6}$ if there are 3 million non-shared SNPs.
945 This burden is much higher than the usually observed somatic mutation rates.

946

947 First, we analysed how many SNPs across 2,504 individuals from the 1000 Genomes Project
948 would remain after filtering with our common SNPs mask ($n=26,111,286$; **Methods**). Our
949 results show that on average 55,685 SNPs would remain unfiltered for a given contaminant
950 individual. Hence, for 1% contamination, filtering of common SNPs would reduce $Burden_{SNP}$
951 from 5×10^{-6} to 9×10^{-8} SNPs/bp. We note that the number of unfiltered SNPs varies largely
952 across continental groups, with averages of 25,666 and 82,765 per individual in Europe and
953 South Asia, respectively (**Supplementary Note 6**).

954

955 To estimate the extent of contamination we rely on VerifyBamID2⁴⁶, which we evaluated
956 simulating contamination fractions below 1%, for both bams sequenced with standard
957 methods and with the NanoSeq protocol (**Extended Dat Fig. 4e, f, Supplementary Note 6**).
958 To obtain more stable estimates we increased the number of markers from 100K to 500K, by
959 randomly choosing additional SNPs with MAF > 0.05 from the 1000 Genomes Project
960 20130502 release.

961

962 ***In silico* decontamination**

963

964 We detected that some libraries were contaminated with DNA from other analysed samples.
965 In cases where the contaminant can be identified, it is possible to remove the mutation calls
966 corresponding to contaminant SNPs by using the corresponding BAM files. This simple
967 approach proved useful to clean contaminated substitution calls and resulting mutation
968 burden corrections were in line with VerifyBamId contamination estimates. That is, mutation
969 burdens of non-contaminated samples remained unaltered after in silico decontamination,
970 whereas the mutation burdens of contaminated samples decreased proportionally to the
971 estimated contamination level.

972

973 This approach was applied to two plates where some samples showed signs of contamination,
 974 including neurons, colonic crypts and smooth muscle samples. Mutation calls occurring at
 975 SNP sites in any of the other samples in the plate were removed. To accomplish this, we
 976 required that each mutation was supported by fewer than 10 base calls across the matched
 977 normals of potential contaminants and that the maximum support from any one matched
 978 normal was lower than 3 reads. All the samples from plates showing evidence of
 979 contamination are considered as potential contaminants. Thresholds to remove contaminant
 980 calls were found empirically for the data at hand and should be adjusted when larger panels
 981 of matched normals or very high coverage samples are analysed.

982
 983 Indels were not analysed for nine samples with signs of contamination as we did not
 984 implement a decontamination pipeline for indels (**Supplementary Table 4**).

985

986 **Correction of mutation burden and trinucleotide substitution profiles**

987

988 Each library preparation method has its own fragmentation and amplification biases and
 989 captures a different subset of the total genome. For instance, amplification biases during
 990 library preparation often lead to lower coverage in GC-rich genomic regions⁴⁷. Since
 991 substitution rates show strong trinucleotide context dependence, taking into consideration
 992 differences in sequence composition can be important when comparing mutation burdens and
 993 substitution profiles between sequencing protocols. Biases can be particularly noticeable with
 994 NanoSeq restriction enzyme libraries, where trinucleotides overlapping the restriction
 995 enzyme site (TGCA in the case of HpyCH4V) are depleted when read ends are filtered. There
 996 are 32 different trinucleotides where the central nucleotide is a pyrimidine. Let t denote the
 997 count of a given trinucleotide of type $i = 1..32$. The frequency of each trinucleotide is
 998 calculated separately for the genome f_i^g and for the NanoSeq experiment (weighted by the
 999 coverage at each site) f_i^e where:

1000

$$f_i = \frac{t_i}{\sum_{i=1}^{32} t_i}$$

1001

1002 The ratio of genomic to experimental frequencies for a given trinucleotide is:

1003

$$r_i = \frac{f_i^g}{f_i^e}$$

1004

1005 There are six classes of substitution where the mutated base is a pyrimidine (C>A, C>G,
 1006 C>T, T>A, T>C, T>G), and for each trinucleotide context there are three possible
 1007 substitutions. Each trinucleotide-substitution count (e.g. ATG>C, where T>C) is corrected by
 1008 the ratio of genomic to experimental frequencies for the corresponding trinucleotide (ATG).
 1009 For instance, let $s_{ATG>C}$ denote the count of substitution $T>C$ in trinucleotide context ATG ,
 1010 the substitution count is corrected as follows:

1011

$$s'_{ATG>C} = s_{ATG>C} r_{ATG}$$

1012

1013 This correction is applied to each of the 96 possible trinucleotide substitutions (h). The
 1014 corrected substitution counts provide a substitution profile projected onto the human genome,
 1015 and are also used to calculate the corrected mutation burden:

1016

$$\beta' = \frac{\sum_{h=1}^{96} s'_i}{\sum_{i=1}^{32} t_i}$$

1017

1018

1019 **Correction of NanoSeq mutation burden in cord blood by accounting for missed early** 1020 **embryonic mutations**

1021

1022 Given their low burden, a substantial fraction of the mutation burden in cord blood HSC/MPP
1023 colonies is attributable to early embryonic mutations shared by multiple colonies. In the
1024 NanoSeq bioinformatic protocol, mutations with a VAF higher than 0.01 in the matched
1025 normal are considered germline SNPs and are filtered out from further analysis. Not
1026 accounting for the loss of early embryonic mutations can have a measurable impact on
1027 burden estimates in cord blood. Taking advantage of the availability of multiple HSC/MPP
1028 colonies per donor, we could quantify the loss of embryonic variants and correct the burden
1029 estimate accordingly. For each of the 50 blood colonies we estimated the global VAF of each
1030 mutation in the remaining 49 colonies. This was done for the two neonatal donors. We
1031 determined that 24% of all the mutations called had a global VAF higher than 0.01. Since a
1032 similar fraction of mutations would be missed by NanoSeq, we multiplied the NanoSeq
1033 estimated burden by a factor of 1.32, i.e. $1/(1-0.24)$. A similar correction is not possible for
1034 the sperm burden estimates, as we lack single-cell level information for sperm, but a modest
1035 underestimation of the mutation burden due to missed embryonic variants is plausible.
1036

1037 **Mutation calling in clonal samples sequenced with standard protocols**

1038

1039 Mutation calls for HSPC colonies from donor PD43976_59yo were obtained from Lee-Six *et*
1040 *al.* 2018²⁶. Mutation calls from standard whole-genome sequencing for the colonic crypts
1041 processed in Lee-Six *et al.* 2019⁶ were obtained from Olafsson *et al.*³⁹. Indel mutation calls
1042 for a bladder tumour sample (**Extended Data Fig. 5**) were obtained from Lawson *et al.*³⁴.
1043 Indel calls for POLE and POLD1 mutants were obtained from Robinson *et al.*⁴⁸ (**Extended**
1044 **Data Fig. 5**).

1045

1046 For the HSC/MPP blood colonies sequenced in the present study, in-house pipelines were
1047 used to run CaVEMan and Pindel against an unmatched synthetic normal genome^{49,50}.
1048 Another bespoke algorithm (cgpVAF) was then used to generate matrices of variant and
1049 normal reads at all sites that had a detected variant in any sample from a given individual.
1050 Up-to-date versions of these algorithms are available from the Sanger Institute's Cancer IT
1051 GitHub repository (<https://github.com/cancerit>).

1052

1053 Filtering strategies detailed below were then used to remove germline variants, technical
1054 artefacts and mutations that had arisen during culture in vitro. (a) A custom filter was used to
1055 remove artefacts associated with the 'low input' library preparation used, including those due
1056 to cruciform DNA structures. (b) A binomial filtering strategy was used to remove variants
1057 with aggregated count distributions consistent with germline single nucleotide
1058 polymorphisms. (c) A beta-binomial filter was used to remove low-frequency artefacts, i.e.
1059 variants present at low frequencies across samples in a way not consistent with the sample-to-
1060 sample variation expected for acquired somatic mutations. (d) Sites with a mean depth below
1061 8 and over 40 were removed. (e) thresholds were used to filter out in vitro variants from the
1062 remaining mutations using a bespoke script. These were set to require a minimum variant
1063 read count of 2 or more and a variant allele fraction of 0.2 for autosomes and 0.4 for XY

1064 chromosomes. (f) The final filtering step involved building a phylogenetic tree from the HSC
1065 genomes derived from each individual. Mutations that did not fit the optimal tree structure
1066 were also discarded as likely artefacts.

1067
1068 Tree building was performed using MPBoot, which is a maximum parsimony tree
1069 approximation method⁵¹. Variants were genotyped as ‘present’ in a sample if 2 or more
1070 variant reads supported the variant. Variants were genotyped as ‘absent’ in a sample if 0
1071 variant reads were present at a given site and depth at that site was 6 or more. Sites that did
1072 not fall into either of the above categories were marked as ‘unknown’. Mutations were
1073 assigned back to the tree using an R package (tree_mut), which uses a maximum likelihood
1074 approach and the original count data to assign each mutation to a branch in the MPBoot
1075 generated tree.

1076

1077 **Estimation of mutation burden in standard sequencing data**

1078

1079 Using clonal or nearly-clonal samples, we were able to compare NanoSeq to mutation burden
1080 estimates from standard whole-genome sequencing. This includes libraries prepared by laser
1081 microdissection and low-input enzymatic fragmentation³⁸ or sonication, followed by standard
1082 Illumina sequencing and mutation calling using CaVEMan⁴⁹. The mutation calls described in
1083 the previous section were further processed to make burden estimates comparable across
1084 protocols.

1085

1086 To compare NanoSeq burdens to those from standard libraries, we restricted the analysis to
1087 regions of the genome covered by at least 20 reads in the standard libraries, to minimise the
1088 impact of low coverage on mutation calling sensitivity. We also excluded the fraction of the
1089 genome flagged as *non-analysed* by CaVEMan. Given the thorough filtering strategies
1090 applied for NanoSeq, we further restricted the analysed genome to include only sites callable
1091 in NanoSeq. Finally, given that trinucleotide frequencies in the callable genome of standard
1092 libraries differ from the background genomic frequencies, burden estimates were corrected
1093 accordingly. The difference in trinucleotide frequencies was mainly due to extensive filtering
1094 of common SNPs (frequent at CpG) and the partial depletion of trinucleotides overlapping
1095 the restriction site (TGCA). Remarkably, we found that estimates of mutation burden
1096 increased by ~20% in standard sequencing data when applying these corrections, largely due
1097 to reducing the impact of low sensitivity in certain genomic regions, either due to low
1098 coverage or mapping quality problems (**Extended Data Fig. 5a, b, Supplementary Note 7**).

1099

1100 **Bootstrapped cosine similarity**

1101

1102 Cosine similarities are frequently used to compare mutational profiles, although they do not
1103 consider the noise introduced by the number of mutations available. Small sample sizes can
1104 cause large cosine similarity deviations from their original spectrum. If a query profile (e.g.
1105 NanoSeq result) with n mutations is to be compared to a reference profile, we can estimate
1106 the impact of small sample sizes by bootstrapping. From the reference profile we obtain
1107 1,000 random samples with size n , and then compare each of these samples back to the
1108 reference profile. We can then calculate the cosine similarities between the query and the
1109 reference profiles and compare it to the 95% interval of cosine similarities observed in the
1110 bootstrapped samples.

1111

1112 **Mutational signature analysis**

1113

1114 Mutational signatures of single-base substitutions in their trinucleotide sequence context were
1115 inferred from sets of somatic mutation counts using the sigfit (v2.0) package for R⁵². *De*
1116 *novo* signature extraction was performed for a range of numbers of signatures ($N = 2, \dots, 8$),
1117 using counts of mutations grouped per tissue type (cord blood, adult blood, granulocytes,
1118 colonic crypts, smooth muscle or neurons), and sequencing method (NanoSeq or standard
1119 sequencing). To account for differences in sequence composition across samples, NanoSeq
1120 mutation counts were corrected as described in a previous section (see **Correction of**
1121 **mutation burden and trinucleotide substitution profiles**). To avoid an excessive influence
1122 of tissue types more highly represented in our dataset, mutation counts were randomly
1123 downsampled to a maximum of 2,000 mutations from each tissue type. Samples with
1124 evidence of sporadic mutational processes, such as APOBEC or colibactin were removed
1125 from the dataset. This excluded urothelium, a bladder tumour sample and colonic crypts from
1126 one donor affected by colibactin (PD37449, $n = 3$). The best-supported number of signatures
1127 on the basis of overall goodness-of-fit, as reported by the ‘extract_signatures’ function in
1128 sigfit, was $N = 3$. The three extracted signatures (**Fig. 3e**) were subsequently fitted to the
1129 counts of mutations per sample (using the ‘fit_signatures’ function in sigfit) to infer the
1130 exposure of each signature in each sample.

1131

1132 Mutational signature analysis was also applied to publicly-available single-nuclei mutation
1133 data from neurons¹³. Three signatures closely matching those shown in the original
1134 publication were extracted using the *extract_signatures* function in sigfit, with parameters
1135 nsignatures=3, seed=1469 and iter=10000.

1136

1137 **Linear regression modelling**

1138

1139 Linear regressions were used to estimate the numbers of mutations accumulated per year, to
1140 test whether mutations associated with a given signature increased with age, or to test the
1141 effects of disease status or organ of origin on mutation burdens.

1142

1143 For neurons and urothelium, with only one sample per donor, we used simple multiple linear
1144 regressions (**Supplementary Table 8**), while for the remaining cell types with multiple
1145 samples per donor (smooth muscle, colonic crypts, blood colonies, granulocytes and sperm)
1146 we used linear mixed-effect models, using donor as a random effect.

1147

1148 For simplicity, in the comparison of substitution and indel rates per year across all cell types
1149 shown in **Fig. 3o,p**, we used regression models without a free intercept, after verifying that
1150 the estimated intercepts were not significantly different from zero. All the regression models,
1151 with and without intercepts, and their parameter estimates are summarised in **Supplementary**
1152 **Table 8**.

1153

1154 To test for the significance of a given fixed effect (such as organ of origin), we used the
1155 anova R function, comparing the null model without the fixed effect and the alternative
1156 model with the fixed effect (**Supplementary Table 8**). Confidence intervals for linear mixed-
1157 effects models at different ages were calculated using parametric bootstrapping and 1,000
1158 replicates, as implemented in the ‘predict’ method in bootpredictlme4 R package.

1159

1160 All linear regression and statistical tests were conducted in R using packages: lm, lmer, afex,
1161 bootpredictlme4, and lmerTest.

1162

1163 **Data Availability**

1164

1165 Information on data availability for all samples is available in **Supplementary Table 1**.
1166 NanoSeq sequencing data has been deposited in EGA under accession number
1167 EGAD00001006459. Sperm samples are available under EGAC0000100027. Standard
1168 sequencing data has been deposited in EGA under accession number EGAD00001006595.
1169 For samples publicly available, references to the original sources are provided in
1170 **Supplementary Table 1**. Substitution and indel rates are available in **Supplementary Table**
1171 **4**. Substitution and indel calls for samples sequenced with NanoSeq are available in
1172 **Supplementary Tables 5** and **6**. Trinucleotide substitution profiles are available in
1173 **Supplementary Table 7**.

1174

1175 **Code Availability**

1176

1177 The bioinformatic pipeline to process NanoSeq sequencing data includes all steps from
1178 processing sequencing data, mapping, calling mutations and calculating corrected burden
1179 estimates and substitution profiles. This code is available from
1180 <https://github.com/cancerit/NanoSeq>. Pipelines to call indels, perform signature extraction
1181 and signature fitting with sigfit, simulate the efficiency of the NanoSeq protocol, calculate
1182 mutation burden in specific genomic regions, and to reproduce most of the main plots are
1183 available from https://github.com/fa8sanger/NanoSeq_Paper_Code. Analyses in R were done
1184 with R v3.3 and v3.6. R libraries used include: GenomicRanges⁵³ (v1.38.0), Rsamtools
1185 (v2.2.3), MASS (v7.3-51.5), sigfit⁵² (v2.0), readxl (v1.3.1), deconstructSigs (v1.8.0), lsa
1186 (v0.73.2), deepSNV⁵⁴ (v1.32.0), lme4 (v1.1-26), afex (v0.28-1), lmerTest (v3.1-3),
1187 bootpredictlme4 (v0.1), and Biostrings (v2.54.0). Our pipeline makes use of samtools⁵⁵ v1.9,
1188 bcftools⁵⁶ v1.9, bwa v0.7.5a-r405, and bedtools⁵⁷ v2.29.0. We also used the following
1189 programs: CaVeMan (v 2020), Pindel (v 2020), and MPBoot 1.1.0.

1190

1191

1192 **Method references**

1193

- 1194 38 Ellis, P. *et al.* Reliable detection of somatic mutations in solid tissues by laser-capture
1195 microdissection and low-input DNA sequencing. *Nat Protoc*, doi:10.1038/s41596-
1196 020-00437-6 (2020).
- 1197 39 Olafsson, S. *et al.* Somatic Evolution in Non-neoplastic IBD-Affected Colon. *Cell*
1198 **182**, 672-684.e611, doi:10.1016/j.cell.2020.06.036 (2020).
- 1199 40 Krishnaswami, S. R. *et al.* Using single nuclei for RNA-seq to capture the
1200 transcriptome of postmortem neurons. *Nat Protoc* **11**, 499-524,
1201 doi:10.1038/nprot.2016.015 (2016).
- 1202 41 Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-
1203 MEM. *arXiv: Genomics* (2013).
- 1204 42 Tischler, G. & Leonard, S. biobambam: tools for read pair collation based algorithms
1205 on BAM files. *Source Code Biol Med* **9**, 13-13, doi:10.1186/1751-0473-9-13 (2014).
- 1206 43 Gerstung, M. *et al.* Reliable detection of subclonal single-nucleotide variants in
1207 tumour cell populations. *Nat Commun* **3**, 811, doi:10.1038/ncomms1814 (2012).
- 1208 44 Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation
1209 in 141,456 humans. *Nature* **581**, 434-443, doi:10.1038/s41586-020-2308-7 (2020).
- 1210 45 Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74,
1211 doi:10.1038/nature15393 (2015).

1212 46 Zhang, F. *et al.* Ancestry-agnostic estimation of DNA sample contamination from
1213 sequence reads. *Genome Res* **30**, 185-194, doi:10.1101/gr.246934.118 (2020).

1214 47 Benjamini, Y. & Speed, T. P. Summarizing and correcting the GC content bias in
1215 high-throughput sequencing. *Nucleic Acids Res* **40**, e72, doi:10.1093/nar/gks001
1216 (2012).

1217 48 Robinson, P. S. *et al.* Elevated somatic mutation burdens in normal human cells due
1218 to defective DNA polymerases. *bioRxiv*, 2020.2006.2023.167668,
1219 doi:10.1101/2020.06.23.167668 (2020).

1220 49 Jones, D. *et al.* cgpCaVEManWrapper: Simple Execution of CaVEMan in Order to
1221 Detect Somatic Single Nucleotide Variants in NGS Data. *Curr Protoc Bioinformatics*
1222 **56**, 15.10.11-15.10.18, doi:10.1002/cpbi.20 (2016).

1223 50 Raine, K. M. *et al.* cgpPindel: Identifying Somatic Acquired Insertion and
1224 Deletion Events from Paired End Sequencing. *Curr Protoc Bioinformatics* **52**,
1225 15.17.11-15.17.12, doi:10.1002/0471250953.bi1507s52 (2015).

1226 51 Hoang, D. T. *et al.* MPBoot: fast phylogenetic maximum parsimony tree inference
1227 and bootstrap approximation. *BMC Evol Biol* **18**, 11, doi:10.1186/s12862-018-1131-3
1228 (2018).

1229 52 Gori, K. & Baez-Ortega, A. sigfit: flexible Bayesian inference of mutational
1230 signatures. *bioRxiv*, 372896, doi:10.1101/372896 (2020).

1231 53 Lawrence, M. *et al.* Software for computing and annotating genomic ranges. *PLoS*
1232 *Comput Biol* **9**, e1003118, doi:10.1371/journal.pcbi.1003118 (2013).

1233 54 Gerstung, M., Papaemmanuil, E. & Campbell, P. J. Subclonal variant calling with
1234 multiple samples and prior knowledge. *Bioinformatics* **30**, 1198-1204,
1235 doi:10.1093/bioinformatics/btt750 (2014).

1236 55 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
1237 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).

1238 56 Li, H. A statistical framework for SNP calling, mutation discovery, association
1239 mapping and population genetical parameter estimation from sequencing data.
1240 *Bioinformatics* **27**, 2987-2993, doi:10.1093/bioinformatics/btr509 (2011).

1241 57 Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing
1242 genomic features. *Bioinformatics* **26**, 841-842, doi:10.1093/bioinformatics/btq033
1243 (2010).

1244 58 Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature*
1245 **518**, 317-330, doi:10.1038/nature14248 (2015).

1246

1247 **Acknowledgements**

1248

1249 We thank Liz Anderson, Kirsty Roberts, Calli Latimer, Quan Lin, the CGP-lab, Rocio
1250 Vicario, Frederic Geissmann, Nicos Angelopoulos, German Tischler, Tristram Bellerby,
1251 Maria Abascal and Krishnaa Chatterjee for assistance in the development of NanoSeq or with
1252 this manuscript.

1253

1254 We are grateful to the live donors and the families of the deceased transplant organ donors.
1255 This research was supported by the Cambridge NIHR BRC Cell Phenotyping Hub. We
1256 gratefully acknowledge the participation of all NIHR BioResource Centre
1257 Cambridge volunteers, and thank the NIHR BioResource Centre Cambridge and staff for
1258 their contribution. We thank the National Institute for Health Research and NHS Blood and
1259 Transplant. The views expressed are those of the author(s) and not necessarily those of the
1260 NHS, the NIHR or the Department of Health & Social Care. We gratefully acknowledge the

1261 Cambridge Blood and Stem Cell Biobank for sample donation and support of this work. We
1262 are grateful to the Cambridge Brain Bank for sample donation. We thank the participants
1263 and local coordinators at the TwinsUK.

1264

1265 **Funding:** I.M. is funded by Cancer Research UK (C57387/A21777) and the Wellcome Trust.
1266 P.J.C. is a Wellcome Trust Senior Clinical Fellow. R.R. is a recipient of a CRUK Career
1267 Development fellowship (C66259/A27114). E.L. is supported by a Wellcome/Royal Society
1268 Sir Henry Dale Fellowship (Grant number 107630/Z/15/Z), the European Hematology
1269 Association, BBSRC and by core funding from Wellcome (Grant number 203151/Z/16/Z)
1270 and MRC to the Wellcome-MRC Cambridge Stem Cell Institute. D.G.K. is supported by a
1271 Bloodwise Bennett Fellowship (15008), the Bill and Melinda Gates Foundation (INV-
1272 002189) and an ERC Starting Grant (ERC-2016-STG-715371). The TwinsUK study was
1273 funded by the Wellcome Trust and European Community's Seventh Framework Programme
1274 (FP7/2007-2013). The TwinsUK study also receives support from the National Institute for
1275 Health Research (NIHR)- funded BioResource, Clinical Research Facility and Biomedical
1276 Research Centre based at Guy's and St Thomas' NHS Foundation Trust in partnership with
1277 King's College London.

1278

1279 **Author Contributions**

1280

1281 R.J.O., F.A., and I.M. conceived the project. I.M., P.J.C., R.R., and M.R.S. supervised the
1282 project. F.A., R.J.O., E.M., and I.M. wrote the manuscript; all authors reviewed and edited
1283 the manuscript. R.J.O. led the development of the protocol with help from F.A., A.R.J.L.,
1284 P.E., S.V.L. and I.M. R.J.O. and F.A. developed the bioinformatics pipeline with help from
1285 R.E.A., S.V.L., and D.J. F.A. led the analysis of the data with help from A.R.J.L., I.M., A.B-
1286 O., Y.W., L.M.R.H., E.J.K., T.H.H.C, M.S.C, and M.G. E.M. performed the HSC/MPP
1287 experiments. L.M.R.H. and A.J.C.R. performed the cell sorting of neuronal nuclei. A.R.J.L.
1288 and A.C. performed laser microdissection. E.M., N.F.O., H.E.M., M.D., D.G.K., E.L.,
1289 K.T.M., K.S.P., K.A., R.R., H.L.S. and S.O collected and processed samples. E.M., E.L.,
1290 M.G. and D.G.K assisted on the interpretation of blood data.

1291

1292 **Competing Interests Declaration**

1293

1294 A patent application on NanoSeq has been filed including R.J.O., F.A. and I.M.

1295

1296 **Additional Information**

1297

1298 **Supplementary information.** The online version contains supplementary material available
1299 at xxx

1300

1301 **Correspondence and requests for materials** should be addressed to R.J.O. or I.M.

1302

1303 **Extended data figures legends**

1304

1305 **Extended Data Figure 1 | Substitution imbalances and impact of A-tailing. a-b,**
1306 Imbalances in the distribution of the six complementary substitutions (e.g. G>T vs C>A)
1307 across read positions in BotSeqS and NanoSeq, respectively. **c,** Origin of G>T over C>A
1308 mutation call imbalances in standard sequencing²². **d,** Origin of imbalances in Duplex
1309 Sequencing / BotSeqS as a result of end repair during library preparation. **e,** Single-strand
1310 consensus calls for pyrimidine (top) and purine (bottom) substitutions for the standard

1311 BotSeqS (left panel) protocol and for NanoSeq with standard and modified A-tailing
1312 protocols (middle and right panels, respectively). For example, C>T changes are shown on
1313 the top, while the complementary G>A changes are shown on bottom. By using ddBTPs
1314 C>A, G>A and T>A errors are reduced, lowering the risk of false positive double-strand
1315 consensus calls.

1316

1317 **Extended Data Figure 2 | BotSeqS errors as a function of read end trimming.** **a**, BotSeqS
1318 estimated burden for the granulocyte sample shown in **Fig. 2c** applying different trimmings
1319 to the 5' ends of reads. Even with extensive trimming we predict at least ~600 artefactual
1320 mutation calls per diploid genome. **b**, Substitution imbalances are observed deep into the
1321 reads and cannot be avoided with read trimming. Imbalances vary from experiment to
1322 experiment, as a consequence of DNA damage on the DNA source or during library
1323 preparation (**Supplementary Note 1**). **c**, Substitution profiles including the reference profile
1324 from single-cell derived blood colonies and three BotSeqS profiles after trimming of 20, 40
1325 and 60 bp from the 5' end of reads (in addition to 15 bp trimming of the 3' end). The text in
1326 the figure indicates the observed and expected cosine similarities (**Methods**) cosine similarity
1327 to the reference profile. C>A and C>G errors in BotSeqS remain after extensive trimming.

1328

1329 **Extended Data Figure 3 | Mung Bean NanoSeq.** **a**, Estimated number of mutations per
1330 cord blood cell. Poisson 95% confidence intervals are shown as lines. The red dotted line
1331 shows the number of mutations per cord blood cell estimated with the restriction enzyme
1332 NanoSeq protocol, with Poisson 95% confidence intervals shown as a red shade. In contrast
1333 to **Fig. 1g**, we did not apply the correction for missing embryonic mutations because here we
1334 are comparing two protocols that are equally affected by this limitation. **b**, Substitution
1335 profiles for the standard end repair protocol (BotSeqS) and for Mung Bean, showing the
1336 cosine similarities with the reference profile (**Fig. 1c**).

1337

1338 **Extended Data Figure 4 | Optimization of duplicate rates, DNA input requirements and**
1339 **estimation of human contamination.** **a**, Relationship between sequencing yield, library
1340 complexity, duplicate rates and efficiency, based on a truncated Poisson model (**Methods**).
1341 From left to right: duplicate rate as a function of the sequencing ratio (sequencing reads /
1342 DNA fragments in the library); efficiency (measured as bases called with duplex
1343 coverage/bases sequenced) as a function of the duplicate rate; and efficiency as a function of
1344 sequencing ratio. **b**, Library yield as fmol per 25 μ l as a function of the amount of input DNA
1345 in ng. **c**, Empirical relationship between the estimated fmol in library (measured by qPCR)
1346 and the number of unique molecules in the library estimated with Picard tools (Lander-
1347 Waterman equation) for our choice of restriction enzyme and fragment size selection (250 -
1348 500 bp). **d**, Empirical relationship between duplicate rates and efficiency of the method,
1349 measured as duplex bases called / number of bases sequenced (i.e. the number of paired-end
1350 reads multiplied by 300). The maximum efficiency (~0.04) is lower than the maximum
1351 analytical expectation (0.12; middle panel in **(a)**) because of the trimming of read ends
1352 (barcodes, restriction sites and 8 bps from each end) and the strict filters that we apply to
1353 consider a site callable. **e**, VerifyBamId contamination estimates for different amounts of
1354 simulated contamination from individuals of different ancestry. **f**, Contamination simulation
1355 using two NanoSeq samples to contaminate each other.

1356

1357 **Extended Data Figure 5 | Correction of standard (CaVEMan-based) mutation burden**
1358 **estimates and validation of NanoSeq indel.** **a**, Comparison of the mutation burden
1359 estimates in regions of the genome with at least 20x coverage (*c*) to the trinucleotide-context-
1360 corrected mutation burdens in the subset of *c* covered by NanoSeq and passing all NanoSeq

1361 filters. **b**, Ratio between the rates shown in panel (a), showing that the corrected burden is
1362 approximately 20% higher than the uncorrected burden; box plots show the interquartile
1363 range, median and 95% confidence interval for the median. **c**, Comparison of indel rates
1364 between cord blood colonies (indels were called with the Pindel algorithm) and granulocytes
1365 from neonates (NanoSeq pipeline), showing Poisson 95% confidence intervals. Given the
1366 sparsity of indel calls in cord blood, data from different colonies (n=100) and granulocytes
1367 (n=2 donors, one of them with 5 replicates) were combined into single point estimates. **d**, The
1368 top two panels show the high similarity between the NanoSeq and Pindel indel profiles for a
1369 bladder tumour; the bottom two profiles show the indel spectra in blood from *POLE* and a
1370 *POLD1* germline mutation carriers, very similar to the reported profiles in Robinson *et al*⁴⁸.
1371

1372 **Extended Data Figure 6 | Haematopoietic stem and progenitor cells and colon histology.**

1373 **a**, Gating strategy for the isolation of HSC/MPPs from a representative bone marrow sample.
1374 Text above plots indicates the population depicted. Text inside the plots indicates the name of
1375 the gates shown in pink. The CD34+/CD38- population is defined as the bottom 20% CD38-
1376 as shown. For all initial samples (BM/PB/CB) the index sorted population is the "HSC pool"
1377 gate. Cell population abundance differed between samples but typically viable cells were 60-
1378 90% of total cells and singlets were 98-99% of viable cells. Live cells were 90-99% of viable
1379 cells and myeloid cells were 15-50% of live cells. CD34+ cells were typically 1-15% of
1380 myeloid cells. **b** and **c**, Colon histology sections showing microbiopsied areas of colonic
1381 epithelium and smooth muscle for donors PD34200 and PD37449, respectively. For donor
1382 PD34200 a single crypt, a pool of six crypts, and two smooth muscle areas were sequenced.
1383 For donor PD37449, the two single crypts and the pool of six crypts were sequenced. The
1384 burden estimates for these microbiopsies are shown in **Fig. 2c** and **3j, k**.
1385

1386 **Extended Data Figure 7 | Neuron nuclei sorting, comparison to single-cell data and**

1387 **accumulation of mutations with age. a**, Gating strategy for the isolation of neuronal nuclei
1388 from frontal cortex. Nuclei were sorted by FACS using an Influx cell sorter (BD Biosciences)
1389 with a 100- μ m nozzle. For each sample an unstained control was used to help determine the
1390 NeuN+ population. The text above each column indicates the population depicted and the
1391 text inside the plots indicates the population of the gates highlighted in black. Sorting results
1392 varied among samples, with 1-60% passing the DAPI gate and, of these, 2-53% passing a
1393 conservative NeuN+ gate. **b**, Substitution profiles for all mutations detected in neurons with
1394 SNP-phased error-corrected single-cell sequencing data in Lodato *et al.*¹³ (top) and with
1395 NanoSeq (middle). In the bottom panel, a signature specific of single-cell sequencing data is
1396 shown (scF signature from Petjak *et al.*¹⁶). **c**, Mutational signatures extracted from Lodato *et*
1397 *al.*¹³, showing their relative contributions in the published dataset. These signatures were
1398 obtained using sigfit (**Methods**) on publicly-available mutation calls and are referred to as
1399 LDA, LDB and LDC. Note the high similarity between the NanoSeq full spectrum for
1400 neurons and LDA (cosine similarity 0.96), and between scF and LDB (cosine similarity
1401 0.97). **d**, Predicted contribution of LDA, LDB and LDC to each of the neurons sequenced in
1402 Lodato *et al.*¹³. **e**, Accumulation of mutations attributed to NanoSeq signatures A, B, and C
1403 with age in healthy donors and in Alzheimer's disease donors. **f**, Accumulation of mutations
1404 attributed to NanoSeq signatures A, B, and C in smooth muscle from bladder and colon.
1405

1406 **Extended Data Figure 8 | Normalised substitution spectra across different genomic**

1407 **regions. a**, Substitution spectra for neurons, granulocytes, smooth muscle and colonic crypts
1408 in chromatin states associated to transcription (states E4 and E5 in ENCODE), and inactive
1409 DNA (E9 and E15). Chromatin states were obtained from ENCODE⁵⁸, using the following
1410 epigenomes: E073 (frontal cortex), E030 (granulocytes), E076 (smooth muscle), and E075

1411 (colonic mucosa). To enable direct comparison of spectra across genomic regions with
1412 different trinucleotide frequencies, the profiles have been normalised to the genomic
1413 trinucleotide frequencies (**Methods**). **b**, Transcriptional strand asymmetries in neurons,
1414 granulocytes, smooth muscle and colonic crypts. **c**, Transcriptional strand asymmetries in
1415 neurons in quartiles of gene expression.

1416

1417 **Extended Data Figure 9 | Additional substitution and indel spectra.** **a**, NanoSeq
1418 mutational spectrum for neurons corrected for trinucleotide frequency in the callable genome.
1419 Unlike the usual representation, which shows unnormalized rates, this representation shows
1420 mutation rates per available trinucleotide. **b**, LDA signature from Lodato *et al.*¹³ normalised
1421 for trinucleotide frequency in the genome also reveals high C>T rates at CpG dinucleotides.
1422 This observation from single-cell data suggests that the high C>T rates at CpG sites in
1423 NanoSeq neuron data (**a**) is not caused by contamination of NeuN+ pools with glia or other
1424 cells. **c**, Indel profiles of granulocytes (top) and of colonic crypts without the colibactin
1425 signature (bottom). **d**, Indel profiles for the 250 most highly expressed genes in PCAWG
1426 liver hepatocellular carcinoma data³¹.

1427

1428 **Extended Data Figure 10 | Smooth muscle.** **a**, Histology of bladder smooth muscle showing
1429 two sections from donor PD40842; only one of the two sections was sequenced with
1430 NanoSeq. **b**, Number of mutations detected with CaVEMan in different smooth muscle
1431 sections processed with our standard microdissection sequencing protocol³⁸. The orange dots
1432 show the expected mutation burdens (with 95% confidence intervals) for these sections based
1433 on the donor age and the regression model shown in **Fig. 3j**. **c**, Distribution of variant allele
1434 frequencies (VAFs) for each of the smooth muscle sections using standard whole-genome
1435 sequencing; box plots show the interquartile range, median, 95% confidence interval for the
1436 median, and outliers as black dots. Box plot notches show the 95% confidence interval for the
1437 median.

1438

1439 **Supplementary tables legends**

1440

1441 **Supplementary Table 1.** Samples used in this study and corresponding data availability

1442 **Supplementary Table 2.** Sequencing yields for NanoSeq/BotSeqS DNA libraries

1443 **Supplementary Table 3.** *In silico* restriction enzyme digestion of the human genome

1444 **Supplementary Table 4.** Substitution and indel rates

1445 **Supplementary Table 5.** Substitution calls (NanoSeq protocol)

1446 **Supplementary Table 6.** Indel calls (NanoSeq protocol)

1447 **Supplementary Table 7.** Trinucleotide substitution profiles

1448 **Supplementary Table 8.** Linear regression models





